

## Chapter 5

# Transforms and Filters for Stochastic Processes

In this chapter, we consider the optimal processing of random signals. We start with transforms that have optimal approximation properties, in the least-squares sense, for continuous and discrete-time signals, respectively. Then we discuss the relationships between discrete transforms, optimal linear estimators, and optimal linear filters.

### 5.1 The Continuous-Time Karhunen–Loève Transform

Among all linear transforms, the *Karhunen–Loève transform* (KLT) is the one which best approximates a stochastic process in the least squares sense. Furthermore, the KLT is a signal expansion with uncorrelated coefficients. These properties make it interesting for many signal processing applications such as coding and pattern recognition. The transform can be formulated for continuous-time and discrete-time processes. In this section, we sketch the continuous-time case [81], [149]. The discrete-time case will be discussed in the next section in greater detail.

Consider a real-valued continuous-time random process  $x(t)$ ,  $a \leq t \leq b$ .

We may not assume that every sample function of the random process lies in  $L_2(a, b)$  and can be represented exactly via a series expansion. Therefore, a weaker condition is formulated, which states that we are looking for a series expansion that represents the stochastic process in the mean:<sup>1</sup>

$$x(t) = \text{l.i.m.}_{N \rightarrow \infty} \sum_{i=1}^N x_i \varphi_i(t) \quad (5.1)$$

The “unknown” orthonormal basis  $\{\varphi_i(t); i = 1, 2, \dots\}$  has to be derived from the properties of the stochastic process. For this, we require that the coefficients

$$x_i = \langle x, \varphi_i \rangle = \int_a^b x(t) \varphi_i(t) dt \quad (5.2)$$

of the series expansion are uncorrelated. This can be expressed as

$$\begin{aligned} E\{x_i x_j\} &= E\{\langle x, \varphi_i \rangle \langle x, \varphi_j \rangle\} \\ &= E\left\{\left(\int_a^b x(t) \varphi_i(t) dt\right) \cdot \left(\int_a^b x(u) \varphi_j(u) du\right)\right\} \\ &= E\left\{\int_a^b \varphi_i(t) \int_a^b x(t) x(u) \varphi_j(u) du dt\right\} \\ &= \int_a^b \varphi_i(t) \left(\int_a^b E\{x(t) x(u)\} \varphi_j(u) du\right) dt \\ &\stackrel{!}{=} \lambda_j \delta_{ij}. \end{aligned} \quad (5.3)$$

The kernel of the integral representation in (5.3) is the autocorrelation function

$$r_{xx}(t, u) = E\{x(t) x(u)\}. \quad (5.4)$$

We see that (5.3) is satisfied if

$$\lambda_j \delta_{ij} = \int_a^b \varphi_i(t) \left(\int_a^b r_{xx}(t, u) \varphi_j(u) du\right) dt. \quad (5.5)$$

Comparing (5.5) with the orthonormality relation  $\delta_{ij} = \int_a^b \varphi_i(t) \varphi_j(t) dt$ , we realize that

$$\int_a^b r_{xx}(t, u) \varphi_j(u) du = \lambda_j \varphi_j(t) \quad (5.6)$$

---

<sup>1</sup>i.i.m.=limit in the mean[38].

must hold in order to satisfy (5.5). Thus, the solutions  $\varphi_j(t)$ ,  $j = 1, 2, \dots$  of the integral equation (5.6) form the desired orthonormal basis. These functions are also called eigenfunctions of the integral operator in (5.6). The values  $\lambda_j$ ,  $j = 1, 2, \dots$  are the eigenvalues. If the kernel  $r_{xx}(t, u)$  is positive definite, that is, if  $\iint r_{xx}(t, u)x(t)x(u) dt du > 0$  for all  $x(t) \in L_2(a, b)$ , then the eigenfunctions form a complete orthonormal basis for  $L_2(a, b)$ . Further properties and particular solutions of the integral equation are for instance discussed in [149].

Signals can be approximated by carrying out the summation in (5.1) only for  $i = 1, 2, \dots, M$  with finite  $M$ . The mean approximation error produced thereby is the sum of those eigenvalues  $\lambda_j$  whose corresponding eigenfunctions are not used for the representation. Thus, we obtain an approximation with minimal mean square error if those eigenfunctions are used which correspond to the largest eigenvalues.

In practice, solving an integral equation represents a major problem. Therefore the continuous-time KLT is of minor interest with regard to practical applications. However, theoretically, that is, without solving the integral equation, this transform is an enormous help. We can describe stochastic processes by means of uncorrelated coefficients, solve estimation or recognition problems for vectors with uncorrelated components and then interpret the results for the continuous-time case.

## 5.2 The Discrete Karhunen–Loève Transform

We consider a real-valued zero-mean random process

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{x} \in \mathbb{R}^n. \quad (5.7)$$

The restriction to zero-mean processes means no loss of generality, since any process  $\mathbf{z}$  with mean  $\mathbf{m}_z$  can be translated into a zero-mean process  $\mathbf{x}$  by

$$\mathbf{x} = \mathbf{z} - \mathbf{m}_z. \quad (5.8)$$

With an orthonormal basis  $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ , the process can be written as

$$\mathbf{x} = \mathbf{U} \boldsymbol{\alpha}, \quad (5.9)$$

where the representation

$$\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_n]^T \quad (5.10)$$

is given by

$$\boldsymbol{\alpha} = \mathbf{U}^T \mathbf{x}. \quad (5.11)$$

As for the continuous-time case, we derive the KLT by demanding uncorrelated coefficients:

$$E \{ \alpha_i \alpha_j \} = \lambda_j \delta_{ij}, \quad i, j = 1, \dots, n. \quad (5.12)$$

The scalars  $\lambda_j$ ,  $j = 1, \dots, n$  are unknown real numbers with  $\lambda_j \geq 0$ . From (5.9) and (5.12) we obtain

$$E \{ \mathbf{u}_i^T \mathbf{x} \mathbf{x}^T \mathbf{u}_j \} = \lambda_j \delta_{ij}, \quad i, j = 1, \dots, n. \quad (5.13)$$

With

$$\mathbf{R}_{xx} = E \{ \mathbf{x} \mathbf{x}^T \} \quad (5.14)$$

this can be written as

$$\mathbf{u}_i^T \mathbf{R}_{xx} \mathbf{u}_j = \lambda_j \delta_{ij}, \quad i, j = 1, \dots, n. \quad (5.15)$$

We observe that because of  $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$ , equation (5.15) is satisfied if the vectors  $\mathbf{u}_j$ ,  $j = 1, \dots, n$  are solutions to the eigenvalue problem

$$\mathbf{R}_{xx} \mathbf{u}_j = \lambda_j \mathbf{u}_j, \quad j = 1, \dots, n. \quad (5.16)$$

Since  $\mathbf{R}_{xx}$  is a covariance matrix, the eigenvalue problem has the following properties:

1. Only real eigenvalues  $\lambda_i$  exist.
2. A covariance matrix is positive definite or positive semidefinite, that is, for all eigenvalues we have  $\lambda_i \geq 0$ .
3. Eigenvectors that belong to different eigenvalues are orthogonal to one another.
4. If multiple eigenvalues occur, their eigenvectors are linearly independent and can be chosen to be orthogonal to one another.

Thus, we see that  $n$  orthogonal eigenvectors always exist. By normalizing the eigenvectors, we obtain the orthonormal basis of the Karhunen–Loève transform.

**Complex-Valued Processes.** For complex-valued processes  $\mathbf{x} \in \mathbb{C}^n$ , condition (5.12) becomes

$$E \{ \alpha_i \alpha_j^* \} = \lambda_j \delta_{ij}, \quad i, j = 1, \dots, n.$$

This yields the eigenvalue problem

$$\mathbf{R}_{xx} \mathbf{u}_j = \lambda_j \mathbf{u}_j, \quad j = 1, \dots, n$$

with the covariance matrix

$$\mathbf{R}_{xx} = E \{ \mathbf{x} \mathbf{x}^H \}.$$

Again, the eigenvalues are real and non-negative. The eigenvectors are orthogonal to one another such that  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$  is unitary.

From the uncorrelatedness of the complex coefficients we cannot conclude that their real and imaginary parts are also uncorrelated; that is,  $E \{ \Re\{\alpha_i\} \Im\{\alpha_j\} \} = 0$ ,  $i, j = 1, \dots, n$  is not implied.

**Best Approximation Property of the KLT.** We henceforth assume that the eigenvalues are sorted such that  $\lambda_1 \geq \dots \geq \lambda_n$ . From (5.12) we get for the variances of the coefficients:

$$E \{ |\alpha_i|^2 \} = \lambda_i, \quad i = 1, \dots, n. \quad (5.17)$$

For the mean-square error of an approximation

$$\hat{\mathbf{x}} = \sum_{i=1}^m \alpha_i \mathbf{u}_i, \quad m < n, \quad (5.18)$$

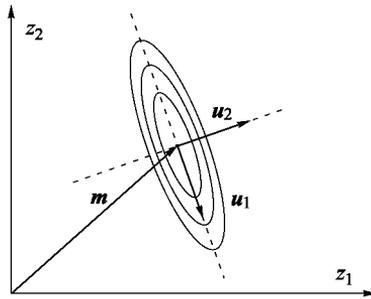
we obtain

$$\begin{aligned} E \{ \|\mathbf{x} - \hat{\mathbf{x}}\|^2 \} &= E \left\{ \left\| \sum_{i=m+1}^n \alpha_i \mathbf{u}_i \right\|^2 \right\} \\ &= \sum_{i=m+1}^n E \{ |\alpha_i|^2 \} \\ &= \sum_{i=m+1}^n \lambda_i. \end{aligned} \quad (5.19)$$

It becomes obvious that an approximation with those eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_m$ , which belong to the largest eigenvalues leads to a minimal error.

In order to show that the KLT indeed yields the smallest possible error among all orthonormal linear transforms, we look at the maximization of  $\sum_{i=1}^m E \{ |\alpha_i|^2 \}$  under the condition  $\|\mathbf{u}_i\| = 1$ . With  $\alpha_i = \mathbf{u}_i^H \mathbf{x}$  this means

$$\sum_{i=1}^m E \{ \mathbf{u}_i^H \mathbf{x} \mathbf{x}^H \mathbf{u}_i \} - \gamma_i \mathbf{u}_i^H \mathbf{u}_i = \sum_{i=1}^m \mathbf{u}_i^H \mathbf{R}_{xx} \mathbf{u}_i - \gamma_i \mathbf{u}_i^H \mathbf{u}_i \stackrel{!}{=} \max, \quad (5.20)$$



**Figure 5.1.** Contour lines of the pdf of a process  $\mathbf{z} = [z_1, z_2]^T$ .

where  $\gamma_i$  are Lagrange multipliers. Setting the gradient to zero yields

$$\mathbf{R}_{xx}\mathbf{u}_i = \gamma_i\mathbf{u}_i, \quad (5.21)$$

which is nothing but the eigenvalue problem (5.16) with  $\gamma_i = \lambda_i$ .

Figure 5.1 gives a geometric interpretation of the properties of the KLT. We see that  $\mathbf{u}_1$  points towards the largest deviation from the center of gravity  $\mathbf{m}$ .

**Minimal Geometric Mean Property of the KLT.** For any positive definite matrix  $\mathbf{X} = X_{ij}$ ,  $i, j = 1, \dots, n$  the following inequality holds [7]:

$$\det \{ \mathbf{X} \} \leq \prod_{k=1}^n X_{kk}. \quad (5.22)$$

Equality is given if  $\mathbf{X}$  is diagonal. Since the KLT leads to a diagonal covariance matrix of the representation, this means that the KLT leads to random variables with a minimal geometric mean of the variances. From this, again, optimal properties in signal coding can be concluded [76].

**The KLT of White Noise Processes.** For the special case that  $\mathbf{R}_{xx}$  is the covariance matrix of a white noise process with

$$\mathbf{R}_{xx} = \sigma^2 \mathbf{I}$$

we have

$$\lambda_1 = \lambda_2 = \dots = \lambda_n = \sigma^2.$$

Thus, the KLT is not unique in this case. Equation (5.19) shows that a white noise process can be optimally approximated with any orthonormal basis.

**Relationships between Covariance Matrices.** In the following we will briefly list some relationships between covariance matrices. With

$$\mathbf{\Lambda} = E \{ \boldsymbol{\alpha} \boldsymbol{\alpha}^H \} = \begin{bmatrix} \lambda_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \lambda_n \end{bmatrix}, \quad (5.23)$$

we can write (5.15) as

$$\mathbf{\Lambda} = \mathbf{U}^H \mathbf{R}_{xx} \mathbf{U}. \quad (5.24)$$

Observing  $\mathbf{U}^H = \mathbf{U}^{-1}$ , we obtain

$$\mathbf{R}_{xx} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H. \quad (5.25)$$

Assuming that all eigenvalues are larger than zero,  $\mathbf{\Lambda}^{-1}$  is given by

$$\mathbf{\Lambda}^{-1} = \begin{bmatrix} \frac{1}{\lambda_1} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \frac{1}{\lambda_n} \end{bmatrix} = \mathbf{U}^H \mathbf{R}_{xx}^{-1} \mathbf{U}. \quad (5.26)$$

Finally, for  $\mathbf{R}_{xx}^{-1}$  we obtain

$$\mathbf{R}_{xx}^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^H. \quad (5.27)$$

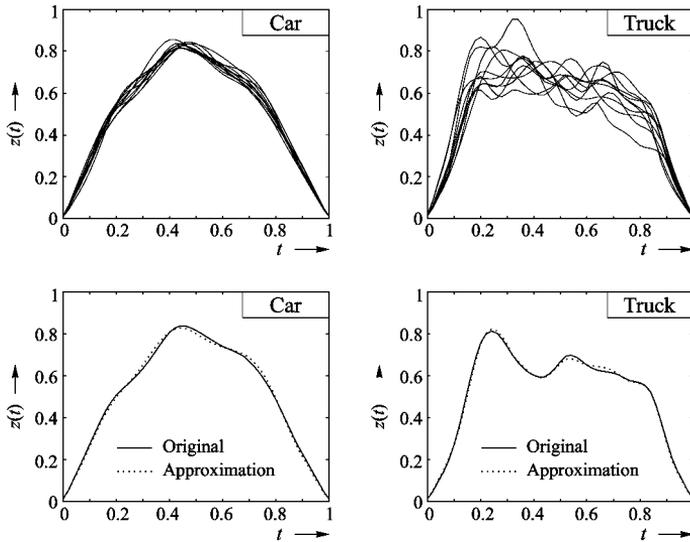
**Application Example.** In pattern recognition it is important to classify signals by means of a few concise features. The signals considered in this example are taken from inductive loops embedded in the pavement of a highway in order to measure the change of inductivity while vehicles pass over them. The goal is to discriminate different types of vehicle (car, truck, bus, etc.). In the following, we will consider the two groups car and truck. After appropriate pre-processing (normalization of speed, length, and amplitude) we obtain the measured signals shown in Figure 5.2, which are typical examples of the two classes. The stochastic processes considered are  $\mathbf{z}_1$  (car) and  $\mathbf{z}_2$  (truck). The realizations are denoted as  ${}^i \mathbf{z}_1$ ,  ${}^i \mathbf{z}_2$ ,  $i = 1 \dots N$ .

In a first step, zero-mean processes are generated:

$$\begin{aligned} \mathbf{x}_1 &= \mathbf{z}_1 - \mathbf{m}_1, \\ \mathbf{x}_2 &= \mathbf{z}_2 - \mathbf{m}_2. \end{aligned} \quad (5.28)$$

The mean values can be estimated by

$$\mathbf{m}_1 = E \{ \mathbf{z}_1 \} \approx \frac{1}{N} \sum_{i=1}^N {}^i \mathbf{z}_1 \quad (5.29)$$



**Figure 5.2.** Examples of sample functions; (a) typical signal contours; (b) two sample functions and their approximations.

and

$$\mathbf{m}_2 = E\{\mathbf{z}_2\} \approx \frac{1}{N} \sum_{i=1}^N \mathbf{z}_2. \quad (5.30)$$

Observing the *a priori* probabilities of the two classes,  $p_1$  and  $p_2$ , a process

$$\mathbf{x} = p_1 \mathbf{x}_1 + p_2 \mathbf{x}_2 \quad (5.31)$$

can be defined. The covariance matrix  $\mathbf{R}_{xx}$  can be estimated as

$$\mathbf{R}_{xx} = E\{\mathbf{x}\mathbf{x}^T\} \approx \frac{p_1}{N+1} \sum_{i=1}^N \mathbf{x}_1 \mathbf{x}_1^T + \frac{p_2}{N+1} \sum_{i=1}^N \mathbf{x}_2 \mathbf{x}_2^T, \quad (5.32)$$

where  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are realizations of the zero-mean processes  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , respectively.

The first ten eigenvalues computed from a training set are:

$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$	$\lambda_8$	$\lambda_9$	$\lambda_{10}$
212923	55460	20559	15790	10230	5262	5036	3139	2551	968

We see that by using only a few eigenvectors a good approximation can be expected. To give an example, Figure 5.2 shows two signals and their

approximations

$$\begin{aligned} {}^1\hat{\mathbf{z}}_1 &= \mathbf{U}\mathbf{U}^T({}^1\mathbf{z}_1 - \mathbf{m}_1) + \mathbf{m}_1, \\ {}^1\hat{\mathbf{z}}_2 &= \mathbf{U}\mathbf{U}^T({}^1\mathbf{z}_2 - \mathbf{m}_2) + \mathbf{m}_2 \end{aligned} \quad (5.33)$$

with the basis  $\{\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3, \mathbf{u}_4\}$ .

In general, the optimality and usefulness of extracted features for discrimination is highly dependent on the algorithm that is used to carry out the discrimination. Thus, the feature extraction method described in this example is not meant to be optimal for all applications. However, it shows how a high proportion of information about a process can be stored within a few features. For more details on classification algorithms and further transforms for feature extraction, see [59, 44, 167, 58].

### 5.3 The KLT of Real-Valued AR(1) Processes

An *autoregressive process* of order  $p$  (AR( $p$ ) process) is generated by exciting a recursive filter of order  $p$  with a zero-mean, stationary white noise process. The filter has the system function

$$H(z) = \frac{1}{1 - \sum_{i=1}^p \rho(i) z^{-i}}, \quad \rho(p) \neq 0. \quad (5.34)$$

Thus, an AR( $p$ ) process  $x(n)$  is described by the difference equation

$$x(n) = w(n) + \sum_{i=1}^p \rho(i) x(n-i), \quad (5.35)$$

where  $w(n)$  is white noise. The AR(1) process with difference equation

$$x(n) = w(n) + \rho x(n-1) \quad (5.36)$$

is often used as a simple model. It is also known as a *first-order Markov process*. From (5.36) we obtain by recursion:

$$x(n) = \sum_{i=0}^{\infty} \rho^i w(n-i). \quad (5.37)$$

For determining the variance of the process  $x(n)$ , we use the properties

$$m_w = E\{w(n)\} = 0 \quad \rightarrow \quad m_x = E\{x(n)\} = 0 \quad (5.38)$$

and

$$r_{ww}(m) = E \{w(n)w(n+m)\} = \sigma^2 \delta_{m0}, \quad (5.39)$$

where  $\delta_{m0}$  is the Kronecker delta. Supposing  $|\rho| < 1$ , we get

$$\begin{aligned} \sigma_x^2 &= E \{|x(n)|^2\} \\ &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \rho^i \rho^j E \{w(n-i)w(n-j)\} \\ &= \sigma^2 \sum_{i=0}^{\infty} \rho^{2i} \\ &= \frac{\sigma^2}{1-\rho^2}. \end{aligned} \quad (5.40)$$

For the autocorrelation sequence we obtain

$$\begin{aligned} r_{xx}(m) &= E \{x(n)x(n+m)\} \\ &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \rho^i \rho^j E \{w(n-i)w(n-j+m)\} \\ &= \sigma^2 \rho^{|m|} \sum_{i=0}^{\infty} \rho^{2i} \\ &= \frac{\sigma^2}{1-\rho^2} \rho^{|m|}. \end{aligned} \quad (5.41)$$

We see that the autocorrelation sequence is infinitely long. However, henceforth only the values  $r_{xx}(-N+1), \dots, r_{xx}(N-1)$  shall be considered. Because of the stationarity of the input process, the covariance matrix of the AR(1) process is a Toeplitz matrix. It is given by

$$\mathbf{R}_{xx} = \frac{\sigma^2}{1-\rho^2} \begin{bmatrix} 1 & \rho & \rho^2 & \dots & \rho^{N-1} \\ \rho & 1 & \rho & & \vdots \\ \rho^2 & \rho & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \rho \\ \rho^{N-1} & \dots & \dots & \rho & 1 \end{bmatrix}. \quad (5.42)$$

The eigenvectors of  $\mathbf{R}_{xx}$  form the basis of the KLT. For real signals and even  $N$ , the eigenvalues  $\lambda_k$ ,  $k = 0, \dots, N-1$  and the eigenvectors were analytically derived by Ray and Driver [123]. The eigenvalues are

$$\lambda_k = \frac{1}{1 - 2\rho \cos(\alpha_k) + \rho^2}, \quad k = 0, \dots, N-1, \quad (5.43)$$

where  $\alpha_k$ ,  $k = 0, \dots, N - 1$  denotes the real positive roots of

$$\tan(N\alpha_k) = -\frac{(1 - \rho^2) \sin(\alpha_k)}{\cos(\alpha_k) - 2\rho + \rho \cos(\alpha_k)}. \quad (5.44)$$

The components of the eigenvectors  $\mathbf{u}_k$ ,  $k = 0, \dots, N - 1$  are given by

$$u_k(n) = \frac{2}{N + \lambda_k} \sin\left(\alpha_k\left(n - \frac{N - 1}{2}\right) + (k + 1)\frac{\pi}{2}\right), \quad n, k = 0, \dots, N - 1. \quad (5.45)$$

## 5.4 Whitening Transforms

In this section we are concerned with the problem of transforming a colored noise process into a white noise process. That is, the coefficients of the representation should not only be uncorrelated (as for the KLT), they should also have the same variance. Such transforms, known as whitening transforms, are mainly applied in signal detection and pattern recognition, because they lead to a convenient process representation with additive white noise.

Let  $\mathbf{n}$  be a process with covariance matrix

$$\mathbf{R}_{nn} = E\{\mathbf{n}\mathbf{n}^H\} \neq \sigma^2 \mathbf{I}. \quad (5.46)$$

We wish to find a linear transform  $\mathbf{T}$  which yields an equivalent process

$$\tilde{\mathbf{n}} = \mathbf{T}\mathbf{n} \quad (5.47)$$

with

$$E\{\tilde{\mathbf{n}}\tilde{\mathbf{n}}^H\} = E\{\mathbf{T}\mathbf{n}\mathbf{n}^H\mathbf{T}^H\} = \mathbf{T}\mathbf{R}_{nn}\mathbf{T}^H = \mathbf{I}. \quad (5.48)$$

We already see that the transform cannot be unique since by multiplying an already computed matrix  $\mathbf{T}$  with an arbitrary unitary matrix, property (5.48) is preserved.

The covariance matrix can be decomposed as follows (KLT):

$$\mathbf{R}_{nn} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^H = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^H\mathbf{U}^H. \quad (5.49)$$

For  $\mathbf{\Lambda}$  and  $\mathbf{\Sigma}$  we have

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_N \end{bmatrix}, \quad \mathbf{\Sigma} = \begin{bmatrix} \sqrt{\lambda_1} & & \\ & \ddots & \\ & & \sqrt{\lambda_N} \end{bmatrix}.$$

Possible transforms are

$$\mathbf{T} = \boldsymbol{\Sigma}^{-1} \mathbf{U}^H \quad (5.50)$$

or

$$\mathbf{T} = \mathbf{U} \boldsymbol{\Sigma}^{-1} \mathbf{U}^H. \quad (5.51)$$

This can easily be verified by substituting (5.50) into (5.48):

$$E \{ \tilde{\mathbf{n}} \tilde{\mathbf{n}}^H \} = \mathbf{T} \mathbf{R}_{nn} \mathbf{T}^H = \boldsymbol{\Sigma}^{-1} \underbrace{\mathbf{U}^H \mathbf{U}}_{\mathbf{I}} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^H \underbrace{\mathbf{U}^H \mathbf{U}}_{\mathbf{I}} \boldsymbol{\Sigma}^{H^{-1}} = \mathbf{I}. \quad (5.52)$$

Alternatively, we can apply the *Cholesky decomposition*

$$\mathbf{R}_{nn} = \mathbf{L} \mathbf{L}^H, \quad (5.53)$$

where  $\mathbf{L}$  is a lower triangular matrix. The whitening transform is

$$\mathbf{T} = \mathbf{L}^{-1}. \quad (5.54)$$

For the covariance matrix we again have

$$E \{ \tilde{\mathbf{n}} \tilde{\mathbf{n}}^H \} = \mathbf{T} \mathbf{R}_{nn} \mathbf{T}^H = \mathbf{L}^{-1} \mathbf{L} \mathbf{L}^H \mathbf{L}^{H^{-1}} = \mathbf{I}. \quad (5.55)$$

In signal analysis, one often encounters signals of the form

$$\mathbf{r} = \mathbf{s} + \mathbf{n}, \quad (5.56)$$

where  $\mathbf{s}$  is a known signal and  $\mathbf{n}$  is an additive colored noise processes. The whitening transforms transfer (5.56) into an equivalent model

$$\tilde{\mathbf{r}} = \tilde{\mathbf{s}} + \tilde{\mathbf{n}} \quad (5.57)$$

with

$$\begin{aligned} \tilde{\mathbf{r}} &= \mathbf{T} \mathbf{r}, \\ \tilde{\mathbf{s}} &= \mathbf{T} \mathbf{s}, \\ \tilde{\mathbf{n}} &= \mathbf{T} \mathbf{n}, \end{aligned} \quad (5.58)$$

where  $\tilde{\mathbf{n}}$  is a white noise process of variance  $\sigma_{\tilde{\mathbf{n}}}^2 = 1$ .

## 5.5 Linear Estimation

In estimation the goal is to determine a set of parameters as precisely as possible from noisy observations. We will focus on the case where the estimators are linear, that is, the estimates for the parameters are computed as linear combinations of the observations. This problem is closely related to the problem of computing the coefficients of a series expansion of a signal, as described in Chapter 3.

Linear methods do not require precise knowledge of the noise statistics; only moments up to the second order are taken into account. Therefore they are optimal only under the linearity constraint, and, in general, non-linear estimators with better properties may be found. However, linear estimators constitute the globally optimal solution as far as Gaussian processes are concerned [149].

### 5.5.1 Least-Squares Estimation

We consider the model

$$\mathbf{r} = \mathbf{S} \mathbf{a} + \mathbf{n}, \quad (5.59)$$

where  $\mathbf{r}$  is our observation,  $\mathbf{a}$  is the parameter vector in question, and  $\mathbf{n}$  is a noise process. Matrix  $\mathbf{S}$  can be understood as a basis matrix that relates the parameters to the clean observation  $\mathbf{S}\mathbf{a}$ .

The requirement to have an unbiased estimate can be written as

$$E\{\hat{\mathbf{a}}(\mathbf{r})|\mathbf{a}\} = \mathbf{a}, \quad (5.60)$$

where  $\mathbf{a}$  is understood as an arbitrary non-random parameter vector. Because of the additive noise, the estimates  $\hat{\mathbf{a}}(\mathbf{r})|\mathbf{a}$  again form a random process.

The linear estimation approach is given by

$$\hat{\mathbf{a}}(\mathbf{r}) = \mathbf{A} \mathbf{r}. \quad (5.61)$$

If we assume zero-mean noise  $\mathbf{n}$ , matrix  $\mathbf{A}$  must satisfy

$$\mathbf{A} \mathbf{S} = \mathbf{I} \quad (5.62)$$

in order to ensure unbiased estimates. This is seen from

$$\begin{aligned}
 E\{\hat{\mathbf{a}}(\mathbf{r})|\mathbf{a}\} &= E\{\mathbf{A}\mathbf{r}|\mathbf{a}\} \\
 &= \mathbf{A}E\{\mathbf{r}|\mathbf{a}\} \\
 &= \mathbf{A}E\{\mathbf{S}\mathbf{a} + \mathbf{n}\} \\
 &= \mathbf{A}\mathbf{S}\mathbf{a} \\
 &\stackrel{!}{=} \mathbf{a}.
 \end{aligned} \tag{5.63}$$

The *generalized least-squares estimator* is derived from the criterion

$$\|\mathbf{r} - \mathbf{S}\boldsymbol{\alpha}\| \Big|_{\boldsymbol{\alpha} = \hat{\mathbf{a}}(\mathbf{r})} \stackrel{!}{=} \min, \tag{5.64}$$

where an arbitrary weighting matrix  $\mathbf{G}$  may be involved in the definition of the inner product that induces the norm in (5.64). Here the observation  $\mathbf{r}$  is considered as a single realization of the stochastic process  $\mathbf{r}$ . Making use of the fact that orthogonal projections yield a minimal approximation error, we get

$$\hat{\mathbf{a}}(\mathbf{r}) = [\mathbf{S}^H \mathbf{G} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{G} \mathbf{r} \tag{5.65}$$

according to (3.95). Assuming that  $[\mathbf{S}^H \mathbf{G} \mathbf{S}]^{-1}$  exists, the requirement (5.65) to have an unbiased estimator is satisfied for arbitrary weighting matrices, as can easily be verified.

If we choose  $\mathbf{G} = \mathbf{I}$ , we speak of a *least-squares estimator*. For weighting matrices  $\mathbf{G} \neq \mathbf{I}$ , we speak of a *generalized least-squares estimator*. However, the approach leaves open the question of how a suitable  $\mathbf{G}$  is found.

### 5.5.2 The Best Linear Unbiased Estimator (BLUE)

As will be shown below, choosing  $\mathbf{G} = \mathbf{R}_{nn}^{-1}$ , where

$$\mathbf{R}_{nn} = E\{\mathbf{n}\mathbf{n}^H\} \tag{5.66}$$

is the correlation matrix of the noise, yields an unbiased estimator with minimal variance. The estimator, which is known as the *best linear unbiased estimator* (BLUE), then is

$$\mathbf{A} = [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{R}_{nn}^{-1}. \tag{5.67}$$

The estimate is given by

$$\hat{\mathbf{a}}(\mathbf{r}) = [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{r}. \tag{5.68}$$

The variances of the individual estimates can be found on the main diagonal of the covariance matrix of the error  $e = \hat{\mathbf{a}}(\mathbf{r}) - \mathbf{a}$ , given by

$$\mathbf{R}_{ee} = [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1}. \quad (5.69)$$

*Proof of (5.69) and the optimality of (5.67).* First, observe that with  $\mathbf{A}\mathbf{S} = \mathbf{I}$  we have

$$\begin{aligned} \hat{\mathbf{a}}(\mathbf{r}) - \mathbf{a} | \mathbf{a} &= \mathbf{A} \mathbf{S} \mathbf{a} + \mathbf{A} \mathbf{n} - \mathbf{a} \\ &= \mathbf{A} \mathbf{n}. \end{aligned} \quad (5.70)$$

Thus,

$$\begin{aligned} \mathbf{R}_{ee} &= \mathbf{A} E \{ \mathbf{n} \mathbf{n}^H \} \mathbf{A}^H \\ &= \mathbf{A} \mathbf{R}_{nn} \mathbf{A}^H \\ &= [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{R}_{nn} \mathbf{R}_{nn}^{-1} \mathbf{S} [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \\ &= [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1}. \end{aligned} \quad (5.71)$$

In order to see whether  $\mathbf{A}$  according to (5.67) is optimal, an estimation

$$\tilde{\mathbf{a}}(\mathbf{r}) = \tilde{\mathbf{A}} \mathbf{r} \quad (5.72)$$

with

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{D} \quad (5.73)$$

will be considered. The unbiasedness constraint requires that

$$\tilde{\mathbf{A}} \mathbf{S} = \mathbf{I}. \quad (5.74)$$

Because of  $\mathbf{A}\mathbf{S} = \mathbf{I}$  this means

$$\mathbf{D}\mathbf{S} = \mathbf{0} \quad (\text{null matrix}). \quad (5.75)$$

For the covariance matrix of the error  $\tilde{e}(\mathbf{r}) = \tilde{\mathbf{a}}(\mathbf{r}) - \mathbf{a}$  we obtain

$$\begin{aligned} \mathbf{R}_{\tilde{e}\tilde{e}} &= \tilde{\mathbf{A}} \mathbf{R}_{nn} \tilde{\mathbf{A}}^H \\ &= [\mathbf{A} + \mathbf{D}] \mathbf{R}_{nn} [\mathbf{A} + \mathbf{D}]^H \\ &= \mathbf{A} \mathbf{R}_{nn} \mathbf{A}^H + \mathbf{A} \mathbf{R}_{nn} \mathbf{D}^H + \mathbf{D} \mathbf{R}_{nn} \mathbf{A}^H + \mathbf{D} \mathbf{R}_{nn} \mathbf{D}^H. \end{aligned} \quad (5.76)$$

With

$$\begin{aligned}
 (\mathbf{A}\mathbf{R}_{nn}\mathbf{D}^H)^H &= \mathbf{D}\mathbf{R}_{nn}\mathbf{A}^H = \mathbf{D}\mathbf{R}_{nn}\mathbf{R}_{nn}^{-1}\mathbf{S}[\mathbf{S}^H\mathbf{R}_{nn}^{-1}\mathbf{S}]^{-1} \\
 &= \underbrace{\mathbf{D}\mathbf{S}}_{\mathbf{0}}[\mathbf{S}^H\mathbf{R}_{nn}^{-1}\mathbf{S}]^{-1} \\
 &= \mathbf{0}
 \end{aligned} \tag{5.77}$$

(5.76) reduces to

$$\mathbf{R}_{\tilde{e}\tilde{e}} = \mathbf{A}\mathbf{R}_{nn}\mathbf{A}^H + \mathbf{D}\mathbf{R}_{nn}\mathbf{D}^H. \tag{5.78}$$

We see that  $\mathbf{R}_{\tilde{e}\tilde{e}}$  is the sum of two non-negative definite expressions so that minimal main diagonal elements of  $\mathbf{R}_{\tilde{e}\tilde{e}}$  are yielded for  $\mathbf{D} = \mathbf{0}$  and thus for  $\mathbf{A}$  according to (5.67).  $\square$

In the case of a white noise process  $\mathbf{n}$ , (5.68) reduces to

$$\hat{\mathbf{a}}(\mathbf{r}) = [\mathbf{S}^H\mathbf{S}]^{-1}\mathbf{S}^H\mathbf{r}. \tag{5.79}$$

Otherwise the weighting with  $\mathbf{G} = \mathbf{R}_{nn}^{-1}$  can be interpreted as an implicit whitening of the noise. This can be seen by using the Cholesky decomposition  $\mathbf{R}_{nn} = \mathbf{L}\mathbf{L}^H$  and by rewriting the model as

$$\tilde{\mathbf{r}} = \tilde{\mathbf{S}}\mathbf{a} + \tilde{\mathbf{n}}, \tag{5.80}$$

where

$$\tilde{\mathbf{r}} = \mathbf{L}^{-1}\mathbf{r}, \quad \tilde{\mathbf{S}} = \mathbf{L}^{-1}\mathbf{S}, \quad \tilde{\mathbf{n}} = \mathbf{L}^{-1}\mathbf{n}. \tag{5.81}$$

The transformed process  $\tilde{\mathbf{n}}$  is a white noise process. The equivalent estimator then is

$$\hat{\mathbf{a}}(\tilde{\mathbf{r}}) = [\tilde{\mathbf{S}}^H\tilde{\mathbf{S}}]^{-1}\tilde{\mathbf{S}}^H\tilde{\mathbf{r}}. \tag{5.82}$$

### 5.5.3 Minimum Mean Square Error Estimation

The advantage of the linear estimators considered in the previous section is their unbiasedness. If we dispense with this property, estimates with smaller mean square error may be found. We will start the discussion on the assumptions

$$E\{\mathbf{r}\} = \mathbf{0}, \quad E\{\mathbf{a}\} = \mathbf{0}. \tag{5.83}$$

Again, the linear estimator is described by a matrix  $\mathbf{A}$ :

$$\hat{\mathbf{a}}(\mathbf{r}) = \mathbf{A}\mathbf{r}. \tag{5.84}$$

Here,  $\mathbf{r}$  is somehow dependent on  $\mathbf{a}$ , but the inner relationship between  $\mathbf{r}$  and  $\mathbf{a}$  need not be known however. The matrix  $\mathbf{A}$  which yields minimal main diagonal elements of the correlation matrix of the estimation error  $\mathbf{e} = \mathbf{a} - \hat{\mathbf{a}}$  is called the *minimum mean square error (MMSE) estimator*.

In order to find the optimal  $\mathbf{A}$ , observe that

$$\begin{aligned} \mathbf{R}_{ee} &= E \left\{ [\hat{\mathbf{a}} - \mathbf{a}] [\hat{\mathbf{a}} - \mathbf{a}]^H \right\} \\ &= E \left\{ \mathbf{a}\mathbf{a}^H \right\} - E \left\{ \hat{\mathbf{a}}\mathbf{a}^H \right\} - E \left\{ \mathbf{a}\hat{\mathbf{a}}^H \right\} + E \left\{ \hat{\mathbf{a}}\hat{\mathbf{a}}^H \right\}. \end{aligned} \quad (5.85)$$

Substituting (5.84) into (5.85) yields

$$\mathbf{R}_{ee} = \mathbf{R}_{aa} - \mathbf{A}\mathbf{R}_{ar} - \mathbf{R}_{ra}\mathbf{A}^H + \mathbf{A}\mathbf{R}_{rr}\mathbf{A}^H \quad (5.86)$$

with

$$\begin{aligned} \mathbf{R}_{aa} &= E \left\{ \mathbf{a}\mathbf{a}^H \right\}, \\ \mathbf{R}_{ar} &= \mathbf{R}_{ra}^H = E \left\{ \mathbf{r}\mathbf{a}^H \right\}, \\ \mathbf{R}_{rr} &= E \left\{ \mathbf{r}\mathbf{r}^H \right\}. \end{aligned} \quad (5.87)$$

Assuming the existence of  $\mathbf{R}_{rr}^{-1}$ , (5.86) can be extended by

$$\mathbf{R}_{ra}\mathbf{R}_{rr}^{-1}\mathbf{R}_{ar} - \mathbf{R}_{ra}\mathbf{R}_{rr}^{-1}\mathbf{R}_{ar}$$

and be re-written as

$$\mathbf{R}_{ee} = [\mathbf{A} - \mathbf{R}_{ra}\mathbf{R}_{rr}^{-1}] \mathbf{R}_{rr} [\mathbf{A}^H - \mathbf{R}_{rr}^{-1}\mathbf{R}_{ar}] - \mathbf{R}_{ra}\mathbf{R}_{rr}^{-1}\mathbf{R}_{ar} + \mathbf{R}_{aa}. \quad (5.88)$$

Clearly,  $\mathbf{R}_{ee}$  has positive diagonal elements. Since only the first term on the right-hand side of (5.88) is dependent on  $\mathbf{A}$ , we have a minimum of the diagonal elements of  $\mathbf{R}_{ee}$  for

$$\mathbf{A} = \mathbf{R}_{ra}\mathbf{R}_{rr}^{-1}. \quad (5.89)$$

The correlation matrix of the estimation error is then given by

$$\mathbf{R}_{ee} = \mathbf{R}_{aa} - \mathbf{R}_{ra}\mathbf{R}_{rr}^{-1}\mathbf{R}_{ar}. \quad (5.90)$$

**Orthogonality Principle.** In Chapter 3 we saw that approximations  $\hat{\mathbf{x}}$  of signals  $\mathbf{x}$  are obtained with minimal error if the error  $\hat{\mathbf{x}} - \mathbf{x}$  is orthogonal to  $\hat{\mathbf{x}}$ . A similar relationship holds between parameter vectors  $\mathbf{a}$  and their MMSE estimates. With  $\mathbf{A}$  according to (5.89) we have

$$\mathbf{R}_{ra} = \mathbf{A} \mathbf{R}_{rr}, \quad \text{i.e.} \quad E \left\{ \mathbf{a}\mathbf{r}^H \right\} = \mathbf{A} E \left\{ \mathbf{r}\mathbf{r}^H \right\}. \quad (5.91)$$

This means that the following orthogonality relations hold:

$$\begin{aligned} E \{ [\hat{\mathbf{a}} - \mathbf{a}] \hat{\mathbf{a}}^H \} &= \mathbf{R}_{\hat{\mathbf{a}}\hat{\mathbf{a}}} - \mathbf{R}_{\hat{\mathbf{a}}\mathbf{a}} \\ &= [\mathbf{A}\mathbf{R}_{\mathbf{r}\mathbf{r}} - \mathbf{R}_{\mathbf{r}\mathbf{a}}] \mathbf{A}^H \\ &= \mathbf{0}. \end{aligned} \quad (5.92)$$

With  $\mathbf{A}\mathbf{r} = \hat{\mathbf{a}}$  the right part of (5.91) can also be written as

$$E \{ \mathbf{a}\mathbf{r}^H \} = E \{ \hat{\mathbf{a}}\mathbf{r}^H \}, \quad (5.93)$$

which yields

$$E \{ [\hat{\mathbf{a}} - \mathbf{a}] \mathbf{r}^H \} = \mathbf{0}. \quad (5.94)$$

The relationship expressed in (5.94) is referred to as the *orthogonality principle*. The orthogonality principle states that we get an MMSE estimate if the error  $\hat{\mathbf{a}}(\mathbf{r}) - \mathbf{a}$  is uncorrelated to all components of the input vector  $\mathbf{r}$  used for computing  $\hat{\mathbf{a}}(\mathbf{r})$ .

**Singular Correlation Matrix.** There are cases where the correlation matrix  $\mathbf{R}_{\mathbf{r}\mathbf{r}}$  becomes singular and the linear estimator cannot be written as

$$\mathbf{A} = \mathbf{R}_{\mathbf{r}\mathbf{a}} \mathbf{R}_{\mathbf{r}\mathbf{r}}^{-1}. \quad (5.95)$$

A more general solution, which involves the replacement of the inverse by the pseudoinverse, is

$$\mathbf{A} = \mathbf{R}_{\mathbf{r}\mathbf{a}} \mathbf{R}_{\mathbf{r}\mathbf{r}}^+. \quad (5.96)$$

In order to show the optimality of (5.96), the estimator

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{D} \quad (5.97)$$

with  $\mathbf{A}$  according to (5.96) and an arbitrary matrix  $\mathbf{D}$  is considered. Using the properties of the pseudoinverse, we derive from (5.97) and (5.86):

$$\begin{aligned} \mathbf{R}_{ee} &= \mathbf{R}_{aa} - \tilde{\mathbf{A}}\mathbf{R}_{ar} - \mathbf{R}_{ra}\tilde{\mathbf{A}}^H + \tilde{\mathbf{A}}\mathbf{R}_{rr}\tilde{\mathbf{A}}^H \\ &= \mathbf{R}_{aa} - \mathbf{R}_{ra}\mathbf{R}_{rr}^+\mathbf{R}_{ar} + \mathbf{D}\mathbf{R}_{rr}^+\mathbf{D}^H. \end{aligned} \quad (5.98)$$

Since  $\mathbf{R}_{rr}^+$  is at least positive semidefinite, we get a minimum of the diagonal elements of  $\mathbf{R}_{ee}$  for  $\mathbf{D} = \mathbf{0}$ , and (5.96) constitutes one of the optimal solutions.

**Additive Uncorrelated Noise.** So far, nothing has been said about possible dependencies between  $\mathbf{a}$  and the noise contained in  $\mathbf{r}$ . Assuming that

$$\mathbf{r} = \mathbf{S}\mathbf{a} + \mathbf{n}, \quad (5.99)$$

where  $\mathbf{n}$  is an additive, uncorrelated noise process, we have

$$\begin{aligned}\mathbf{R}_{ra} &= \mathbf{R}_{ar}^H = \mathbf{R}_{aa} \mathbf{S}^H, \\ \mathbf{R}_{rr} &= \mathbf{S} \mathbf{R}_{aa} \mathbf{S}^H + \mathbf{R}_{nn},\end{aligned}\quad (5.100)$$

and  $\mathbf{A}$  according to (5.89) becomes

$$\mathbf{A} = \mathbf{R}_{aa} \mathbf{S}^H [\mathbf{S} \mathbf{R}_{aa} \mathbf{S}^H + \mathbf{R}_{nn}]^{-1}. \quad (5.101)$$

Alternatively,  $\mathbf{A}$  can be written as

$$\mathbf{A} = [\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{R}_{nn}^{-1}. \quad (5.102)$$

This is verified by equating (5.101) and (5.102), and by multiplying the obtained expression with  $[\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]$  from the left and with  $[\mathbf{S} \mathbf{R}_{aa} \mathbf{S}^H + \mathbf{R}_{nn}]$  from the right, respectively:

$$[\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}] \mathbf{R}_{aa} \mathbf{S}^H = \mathbf{S}^H \mathbf{R}_{nn}^{-1} [\mathbf{S} \mathbf{R}_{aa} \mathbf{S}^H + \mathbf{R}_{nn}].$$

The equality of both sides is easily seen. The matrices to be inverted in (5.102), except  $\mathbf{R}_{nn}$ , typically have a much smaller dimension than those in (5.101). If the noise is white,  $\mathbf{R}_{nn}^{-1}$  can be immediately stated, and (5.102) is advantageous in terms of computational cost.

For  $\mathbf{R}_{ee}$  we get from (5.89), (5.90), (5.100) and (5.102):

$$\begin{aligned}\mathbf{R}_{ee} &= \mathbf{R}_{aa} - \mathbf{A} \mathbf{R}_{ar} \\ &= \mathbf{R}_{aa} - [\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S} \mathbf{R}_{aa}.\end{aligned}\quad (5.103)$$

Multiplying (5.103) with  $[\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]$  from the left yields

$$\begin{aligned}[\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}] \mathbf{R}_{ee} &= [\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}] \mathbf{R}_{aa} - \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S} \mathbf{R}_{aa} \\ &= \mathbf{I},\end{aligned}\quad (5.104)$$

so that the following expression is finally obtained:

$$\mathbf{R}_{ee} = [\mathbf{R}_{aa}^{-1} + \mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1}. \quad (5.105)$$

**Equivalent Estimation Problems.** We partition  $\mathbf{A}$  and  $\mathbf{a}$  into

$$\mathbf{a} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix}, \quad (5.106)$$

such that

$$\begin{bmatrix} \hat{\mathbf{a}}_1(\mathbf{r}) \\ \hat{\mathbf{a}}_2(\mathbf{r}) \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix} \mathbf{r}. \quad (5.107)$$

If we assume that the processes  $\mathbf{a}_1$ ,  $\mathbf{a}_2$  and  $\mathbf{n}$  are independent of one another, the covariance matrix  $\mathbf{R}_{aa}$  and its inverse  $\mathbf{R}_{aa}^{-1}$  have the form

$$\mathbf{R}_{aa} = \begin{bmatrix} \mathbf{R}_{a_1 a_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_{a_2 a_2} \end{bmatrix}, \quad \mathbf{R}_{aa}^{-1} = \begin{bmatrix} \mathbf{R}_{a_1 a_1}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_{a_2 a_2}^{-1} \end{bmatrix}, \quad (5.108)$$

and  $\mathbf{A}$  according to (5.102) can be written as

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \mathbf{S}_1 + \mathbf{R}_{a_1 a_1}^{-1} & \mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \mathbf{S}_2 \\ \mathbf{S}_2^H \mathbf{R}_{nn}^{-1} \mathbf{S}_1 & \mathbf{S}_2^H \mathbf{R}_{nn}^{-1} \mathbf{S}_2 + \mathbf{R}_{a_2 a_2}^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \\ \mathbf{S}_2^H \mathbf{R}_{nn}^{-1} \end{bmatrix}, \quad (5.109)$$

where  $\mathbf{S} = [\mathbf{S}_1, \mathbf{S}_2]$ . Applying the matrix equation

$$\begin{bmatrix} \mathcal{E} & \mathcal{F} \\ \mathcal{G} & \mathcal{H} \end{bmatrix}^{-1} = \begin{bmatrix} \mathcal{E}^{-1} + \mathcal{E}^{-1} \mathcal{F} \mathcal{D}^{-1} \mathcal{G} \mathcal{E}^{-1} & -\mathcal{E}^{-1} \mathcal{F} \mathcal{D}^{-1} \\ -\mathcal{D}^{-1} \mathcal{G} \mathcal{E}^{-1} & \mathcal{D}^{-1} \end{bmatrix} \quad (5.110)$$

$$\mathcal{D} = \mathcal{H} - \mathcal{G} \mathcal{E}^{-1} \mathcal{F}$$

yields

$$\mathbf{A}_1 = [\mathbf{S}_1^H \mathbf{R}_{n_1 n_1}^{-1} \mathbf{S}_1 + \mathbf{R}_{a_1 a_1}^{-1}]^{-1} \mathbf{S}_1^H \mathbf{R}_{n_1 n_1}^{-1}, \quad (5.111)$$

$$\mathbf{A}_2 = [\mathbf{S}_2^H \mathbf{R}_{n_2 n_2}^{-1} \mathbf{S}_2 + \mathbf{R}_{a_2 a_2}^{-1}]^{-1} \mathbf{S}_2^H \mathbf{R}_{n_2 n_2}^{-1} \quad (5.112)$$

with

$$\mathbf{R}_{n_1 n_1} = \mathbf{R}_{nn} + \mathbf{S}_2 \mathbf{R}_{a_2 a_2} \mathbf{S}_2^H, \quad (5.113)$$

$$\mathbf{R}_{n_2 n_2} = \mathbf{R}_{nn} + \mathbf{S}_1 \mathbf{R}_{a_1 a_1} \mathbf{S}_1^H. \quad (5.114)$$

The inverses  $\mathbf{R}_{n_1 n_1}^{-1}$  and  $\mathbf{R}_{n_2 n_2}^{-1}$  can be written as

$$\mathbf{R}_{n_1 n_1}^{-1} = \left[ \mathbf{R}_{nn}^{-1} - \mathbf{R}_{nn}^{-1} \mathbf{S}_2 (\mathbf{S}_2^H \mathbf{R}_{nn}^{-1} \mathbf{S}_2 + \mathbf{R}_{a_2 a_2}^{-1})^{-1} \mathbf{S}_2^H \mathbf{R}_{nn}^{-1} \right], \quad (5.115)$$

$$\mathbf{R}_{n_2 n_2}^{-1} = \left[ \mathbf{R}_{nn}^{-1} - \mathbf{R}_{nn}^{-1} \mathbf{S}_1 (\mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \mathbf{S}_1 + \mathbf{R}_{a_1 a_1}^{-1})^{-1} \mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \right]. \quad (5.116)$$

Equations (5.111) and (5.112) describe estimations of  $\mathbf{a}_1$  and  $\mathbf{a}_2$  in the models

$$\mathbf{r} = \mathbf{S}_1 \mathbf{a}_1 + \mathbf{n}_1, \quad (5.117)$$

$$\mathbf{r} = \mathbf{S}_2 \mathbf{a}_2 + \mathbf{n}_2 \quad (5.118)$$

with

$$\begin{aligned} \mathbf{n}_1 &= \mathbf{S}_2 \mathbf{a}_2 + \mathbf{n}, \\ \mathbf{n}_2 &= \mathbf{S}_1 \mathbf{a}_1 + \mathbf{n}. \end{aligned} \quad (5.119)$$

Thus, each parameter to be estimated can be understood as noise in the estimation of the remaining parameters. An exception is given if  $\mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \mathbf{S}_2 = \mathbf{0}$ , which means that  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are orthogonal to each other with respect to the weighting matrix  $\mathbf{R}_{nn}^{-1}$ . Then we get

$$\mathbf{A}_1 = [\mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \mathbf{S}_1 + \mathbf{R}_{a_1 a_1}^{-1}]^{-1} \mathbf{S}_1^H \mathbf{R}_{nn}^{-1}$$

and

$$\mathbf{R}_{e_1 e_1} = [\mathbf{S}_1^H \mathbf{R}_{nn}^{-1} \mathbf{S}_1 + \mathbf{R}_{a_1 a_1}^{-1}]^{-1},$$

and we observe that the second signal component  $\mathbf{S}_2 \mathbf{a}_2$  has no influence on the estimate.

**Nonzero-Mean Processes.** One could imagine that the precision of linear estimations with respect to nonzero-mean processes  $\mathbf{r}$  and  $\mathbf{a}$  can be increased compared to the solutions above if an additional term taking care of the mean values of the processes is considered. In order to describe this more general case, let us denote the mean of the parameters as

$$\bar{\mathbf{a}} = E\{\mathbf{a}\}. \quad (5.120)$$

The estimate is now written as

$$\hat{\mathbf{a}} = \mathbf{A} \mathbf{r} + \bar{\mathbf{a}} + \mathbf{c}, \quad (5.121)$$

where  $\mathbf{c}$  is yet unknown. Using the shorthand

$$\begin{aligned} \mathbf{b} &= \mathbf{a} - \bar{\mathbf{a}}, \\ \hat{\mathbf{b}} &= \hat{\mathbf{a}} - \bar{\mathbf{a}}, \\ \mathbf{M} &= [\mathbf{c}, \mathbf{A}] \\ \mathbf{x} &= \begin{bmatrix} 1 \\ \mathbf{r} \end{bmatrix}, \end{aligned} \quad (5.122)$$

(5.121) can be rewritten as:

$$\hat{\mathbf{b}} = \mathbf{M} \mathbf{x}. \quad (5.123)$$

The relationship between  $\hat{\mathbf{b}}$  and  $\mathbf{x}$  is linear as usual, so that the optimal  $\mathbf{M}$  can be given according to (5.89):

$$\mathbf{M} = \mathbf{R}_{xb} \mathbf{R}_{xx}^{-1}. \quad (5.124)$$

Now let us express  $\mathbf{R}_{xb}$  and  $\mathbf{R}_{xx}^{-1}$  through correlation matrices of the processes  $\mathbf{a}$  and  $\mathbf{r}$ . From (5.122) and  $E\{\mathbf{b}\} = \mathbf{0}$  we derive

$$\mathbf{R}_{xb} = [\mathbf{0}, \mathbf{R}_{rb}] \quad (5.125)$$

with

$$\begin{aligned} \mathbf{R}_{rb} &= E\{[\mathbf{a} - \bar{\mathbf{a}}] \mathbf{r}^H\} \\ &= E\{[\mathbf{a} - \bar{\mathbf{a}}] [\mathbf{r} - \bar{\mathbf{r}}]^H\}, \end{aligned} \quad (5.126)$$

where

$$\bar{\mathbf{r}} = E\{\mathbf{r}\}. \quad (5.127)$$

$\mathbf{R}_{xx}$  writes

$$\mathbf{R}_{xx} = \begin{bmatrix} 1 & \bar{\mathbf{r}}^H \\ \bar{\mathbf{r}} & \mathbf{R}_{rr} \end{bmatrix}. \quad (5.128)$$

Using (5.110) we obtain

$$\mathbf{R}_{xx}^{-1} = \begin{bmatrix} 1 + \bar{\mathbf{r}}^H [\mathbf{R}_{rr} - \bar{\mathbf{r}} \bar{\mathbf{r}}^H]^{-1} \bar{\mathbf{r}} & -\bar{\mathbf{r}}^H [\mathbf{R}_{rr} - \bar{\mathbf{r}} \bar{\mathbf{r}}^H]^{-1} \\ -[\mathbf{R}_{rr} - \bar{\mathbf{r}} \bar{\mathbf{r}}^H]^{-1} \bar{\mathbf{r}} & [\mathbf{R}_{rr} - \bar{\mathbf{r}} \bar{\mathbf{r}}^H]^{-1} \end{bmatrix}. \quad (5.129)$$

From (5.122) – (5.129) and

$$[\mathbf{R}_{rr} - \bar{\mathbf{r}} \bar{\mathbf{r}}^H] = E\{[\mathbf{r} - \bar{\mathbf{r}}] [\mathbf{r} - \bar{\mathbf{r}}]^H\} \quad (5.130)$$

we finally conclude

$$\hat{\mathbf{a}} = E\{[\mathbf{a} - \bar{\mathbf{a}}] [\mathbf{r} - \bar{\mathbf{r}}]^H\} E\{[\mathbf{r} - \bar{\mathbf{r}}] [\mathbf{r} - \bar{\mathbf{r}}]^H\}^{-1} [\mathbf{r} - \bar{\mathbf{r}}] + \bar{\mathbf{a}}. \quad (5.131)$$

Equation (5.131) can be interpreted as follows: the nonzero-mean processes  $\mathbf{a}$  and  $\mathbf{r}$  are first modified so as to become zero-mean processes  $\mathbf{a} - \bar{\mathbf{a}}$  and  $\mathbf{r} - \bar{\mathbf{r}}$ . For the zero-mean processes the estimation problem can be solved as usual. Subsequently the mean value  $\bar{\mathbf{a}}$  is added in order to obtain the final estimate  $\hat{\mathbf{a}}$ .

**Unbiasedness for Random Parameter Vectors.** So far the parameter vector to be estimated was assumed to be non-random. If we consider  $\mathbf{a}$  to be a random process, various other weaker definitions of unbiasedness are possible.

The straightforward requirement

$$E \{ \hat{\mathbf{a}}(\mathbf{r}) \} = \mathbf{A} E \{ \mathbf{r} \} \stackrel{!}{=} E \{ \mathbf{a} \} \quad (5.132)$$

is meaningless, because it is satisfied for any  $\mathbf{A}$  as far as  $\mathbf{r}$  and  $\mathbf{a}$  are zero-mean.

A useful definition of unbiasedness in the case of random parameters is to consider one of the parameters contained in  $\mathbf{a}$  (e.g.  $a_k$ ) as non-random and to regard all other parameters  $a_1, \dots, a_{k-1}, a_{k+1}$  etc. as random variables:

$$E \{ \hat{a}_k(\mathbf{r}) \mid a_k \} \stackrel{!}{=} a_k. \quad (5.133)$$

In order to obtain an estimator which is unbiased in the sense of (5.133), the equivalences discussed above may be applied. Starting with the model

$$\begin{aligned} \mathbf{r} &= \mathbf{S}\mathbf{a} + \mathbf{n} \\ &= \mathbf{s}_k a_k + \tilde{\mathbf{n}}, \end{aligned} \quad (5.134)$$

in which  $\tilde{\mathbf{n}}$  contains the additive noise  $\mathbf{n}$  and the signal component produced by all random parameters  $a_j$ ,  $j \neq k$ , we can write the unbiased estimate as

$$\hat{a}_k = \mathbf{h}_k^H \mathbf{r} \quad (5.135)$$

with

$$\mathbf{h}_k^H = [\mathbf{s}_k^H \mathbf{R}_{\tilde{\mathbf{n}}\tilde{\mathbf{n}}}^{-1} \mathbf{s}_k]^{-1} \mathbf{s}_k^H \mathbf{R}_{\tilde{\mathbf{n}}\tilde{\mathbf{n}}}^{-1}. \quad (5.136)$$

Then,

$$\mathbf{A} = [\mathbf{h}_1, \mathbf{h}_2, \dots]^H \quad (5.137)$$

is an estimator which is unbiased in the sense of (5.133).

**The Relationship between MMSE Estimation and the BLUE.** If  $\mathbf{R}_{aa} = E \{ \mathbf{a}\mathbf{a}^H \}$  is unknown,  $\mathbf{R}_{aa}^{-1} = \mathbf{0}$  is substituted into (5.102), and we obtain the BLUE (cf. (5.67)):

$$\mathbf{A} = [\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}]^{-1} \mathbf{S}^H \mathbf{R}_{nn}^{-1}. \quad (5.138)$$

In the previous discussion it became obvious that it is possible to obtain unbiased estimates of some of the parameters and to estimate the others with minimum mean square error. This result is of special interest if no unbiased estimator can be stated for all parameters because of a singular matrix  $\mathbf{S}^H \mathbf{R}_{nn}^{-1} \mathbf{S}$ .

## 5.6 Linear Optimal Filters

### 5.6.1 Wiener Filters

We consider the problem depicted in Figure 5.3. By linear filtering of the noisy signal  $r(n) = x(n) + w(n)$  we wish to make  $y(n) = r(n) * h(n)$  as similar as possible to a desired signal  $d(n)$ . The quality criterion used for designing the optimal causal linear filter  $h(n)$  is

$$E \{|e(n)|^2\} = E \{|d(n) - y(n)|^2\} \stackrel{!}{=} \min. \quad (5.139)$$

The solution to this optimization problem can easily be stated by applying the orthogonality principle. Assuming a causal FIR filter  $h(n)$  of length  $p$ , we have

$$y(n) = \sum_{i=0}^{p-1} h(i) r(n-i). \quad (5.140)$$

Thus, according to (5.94), the following orthogonality condition must be satisfied by the optimal filter:

$$E \left\{ \left[ d(n) - \sum_{i=0}^{p-1} h(i) r(n-i) \right] r^*(n-j) \right\} = 0, \quad j = 0, 1, \dots, p-1. \quad (5.141)$$

For stationary processes  $r(n)$  and  $d(n)$  this yields the discrete form of the so-called *Wiener-Hopf equation*:

$$\sum_{i=0}^{p-1} h(i) r_{rr}(j-i) = r_{rd}(j), \quad j = 0, 1, \dots, p-1, \quad (5.142)$$

with

$$\begin{aligned} r_{rr}(m) &= E \{ r^*(n) r(n+m) \}, \\ r_{rd}(m) &= E \{ r^*(n) d(n+m) \}. \end{aligned} \quad (5.143)$$

The optimal filter is found by solving (5.142).

An application example is the estimation of data  $d(n)$  from a noisy observation  $r(n) = \sum_{\ell} c(\ell) d(n-\ell) + w(n)$ , where  $c(n)$  is a channel and  $w(n)$  is noise. By using the optimal filter  $h(n)$  designed according to (5.142) the data is recovered with minimal mean square error.

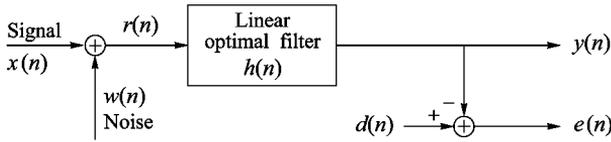


Figure 5.3. Designing linear optimal filters.

**Variance.** For the variance of the error we have

$$\begin{aligned}
 \sigma_e^2 &= E \{ |e(n)|^2 \} \\
 &= \sigma_d^2 - \sum_{i=0}^{p-1} h(i) r_{rd}^*(i) - \sum_{i=0}^{p-1} h^*(i) r_{rd}(i) \\
 &\quad + \sum_{i=0}^{p-1} \sum_{j=0}^{p-1} h(i) h^*(j) r_{rr}(j-i)
 \end{aligned} \tag{5.144}$$

with  $\sigma_d^2 = E \{ |d(n)|^2 \}$ . Substituting the optimal solution (5.142) into (5.144) yields

$$\sigma_{e_{\min}}^2 = \sigma_d^2 - \sum_{i=0}^{p-1} h(i) r_{rd}^*(i). \tag{5.145}$$

**Matrix Notation.** In matrix notation (5.142) is

$$\mathbf{R}_{rr} \mathbf{h} = \mathbf{r}_{rd} \tag{5.146}$$

with

$$\mathbf{h}^T = [h(0), h(1), \dots, h(p-1)], \tag{5.147}$$

$$\mathbf{r}_{rd}^T = [r_{rd}(0), r_{rd}(1), \dots, r_{rd}(p-1)] \tag{5.148}$$

and

$$\mathbf{R}_{rr} = \begin{bmatrix} r_{rr}(0) & r_{rr}(-1) & \dots & r_{rr}(-p+1) \\ r_{rr}(1) & r_{rr}(0) & \dots & r_{rr}(-p+2) \\ \vdots & \vdots & & \vdots \\ r_{rr}(p-1) & r_{rr}(p-1) & \dots & r_{rr}(0) \end{bmatrix}. \tag{5.149}$$

From (5.146) and (5.145) we obtain the following alternative expressions for the minimal variance:

$$\begin{aligned}\sigma_{e_{\min}}^2 &= \sigma_d^2 - \mathbf{r}_{rd}^H \mathbf{h} \\ &= \sigma_d^2 - \mathbf{r}_{rd}^H \mathbf{R}_{rr}^{-1} \mathbf{r}_{rd}.\end{aligned}\tag{5.150}$$

**Special Cases.** The following three cases, where the desired signal is a delayed version of the clean input signal  $x(n)$ , are of special interest:

- (i) Filtering:  $d(n) = x(n)$ .
- (ii) Interpolation:  $d(n) = x(n + D)$ ,  $D < 0$ .
- (iii) Prediction:  $d(n) = x(n + D)$ ,  $D > 0$ . Here the goal is to predict a future value.

For the three cases mentioned above the Wiener–Hopf equation is

$$\sum_{i=0}^{p-1} h(i) r_{rr}(j-i) = r_{rx}(j+D), \quad j = 0, 1, \dots, p-1.\tag{5.151}$$

**Uncorrelated Noise.** If the noise  $w(n)$  is uncorrelated to  $x(n)$ , we have

$$r_{rr}(m) = r_{xx}(m) + r_{ww}(m)\tag{5.152}$$

and

$$r_{rd}(m) = r_{xx}(m+D),\tag{5.153}$$

and from (5.151) we derive

$$\sum_{i=0}^{p-1} h(i) [r_{xx}(j-i) + r_{ww}(j-i)] = r_{xx}(j+D), \quad j = 0, 1, \dots, p-1.\tag{5.154}$$

In matrix notation we get

$$[\mathbf{R}_{xx} + \mathbf{R}_{ww}] \mathbf{h} = \mathbf{r}_{xx}(D)\tag{5.155}$$

with

$$\mathbf{h}^T = [h(0), h(1), \dots, h(p-1)],\tag{5.156}$$

$$\mathbf{r}_{xx}^T(D) = [r_{xx}(D), r_{xx}(D+1), \dots, r_{xx}(D+p-1)]\tag{5.157}$$

and

$$\mathbf{R}_{xx} = \begin{bmatrix} r_{xx}(0) & r_{xx}(-1) & \dots & r_{xx}(-p+1) \\ r_{xx}(1) & r_{xx}(0) & \dots & r_{xx}(-p+2) \\ \vdots & \vdots & & \vdots \\ r_{xx}(p-1) & r_{xx}(p-1) & \dots & r_{xx}(0) \end{bmatrix}. \quad (5.158)$$

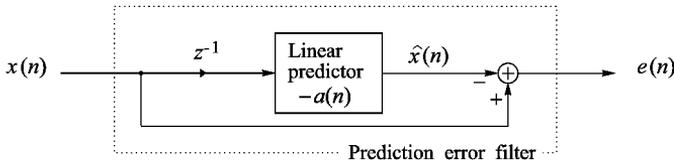
For the correlation matrix  $\mathbf{R}_{ww}$  the corresponding definition holds.

The minimal variance is

$$\begin{aligned} \sigma_{e_{\min}}^2 &= \sigma_d^2 - \mathbf{r}_{xx}^H(D) \mathbf{h} \\ &= \sigma_d^2 - \mathbf{r}_{xx}^H(D) [\mathbf{R}_{xx} + \mathbf{R}_{ww}]^{-1} \mathbf{r}_{xx}(D). \end{aligned} \quad (5.159)$$

### 5.6.2 One-Step Linear Prediction

One-step linear predictors are used in many applications such as speech and image coding (DPCM, ADPCM, LPC,...), in spectral estimation, and in feature extraction for speech recognition. Basically, they may be regarded as a special case of Wiener–Hopf filtering.



**Figure 5.4.** One-step linear prediction.

We consider the system in Figure 5.4. A comparison with Figure 5.3 shows that the optimal predictor can be obtained from the Wiener–Hopf equations for the special case  $D = 1$  with  $d(n) = x(n + 1)$ , while no additive noise is assumed,  $w(n) = 0$ . Note that the filter  $a(n)$  is related to the Wiener–Hopf filter  $h(n)$  as  $a(n) = -h(n - 1)$ . With

$$\hat{x}(n) = - \sum_{i=1}^p a(i) x(n - i), \quad (5.160)$$

where  $p$  is the length of the FIR filter  $a(n)$ , the error becomes

$$\begin{aligned} e(n) &= x(n) - \hat{x}(n) \\ &= x(n) + \sum_{i=1}^p a(i) x(n - i). \end{aligned} \quad (5.161)$$

Minimizing the error with respect to the filter coefficients yields the equations

$$-\sum_{i=1}^p a(i) r_{xx}(j-i) = r_{xx}(j), \quad j = 1, 2, \dots, p, \quad (5.162)$$

which are known as the *normal equations of linear prediction*. In matrix notation they are

$$\begin{bmatrix} r_{xx}(0) & r_{xx}(-1) & \dots & r_{xx}(-p+1) \\ r_{xx}(1) & r_{xx}(0) & \dots & r_{xx}(-p+2) \\ \vdots & \vdots & & \vdots \\ r_{xx}(p-1) & r_{xx}(p-2) & \dots & r_{xx}(0) \end{bmatrix} \begin{bmatrix} a(1) \\ a(2) \\ \vdots \\ a(p) \end{bmatrix} = - \begin{bmatrix} r_{xx}(1) \\ r_{xx}(2) \\ \vdots \\ r_{xx}(p) \end{bmatrix}, \quad (5.163)$$

that is

$$\mathbf{R}_{xx} \mathbf{a} = -\mathbf{r}_{xx}(1) \quad (5.164)$$

with

$$\mathbf{a}^T = [a(1), \dots, a(p)]. \quad (5.165)$$

According to (5.159) we get for the minimal variance:

$$\begin{aligned} \sigma_{e_{\min}}^2 = E\{|e(n)|^2\} &= r_{xx}(0) - \mathbf{r}_{xx}^H(1) \mathbf{R}_{xx}^{-1} \mathbf{r}_{xx}(1) \\ &= r_{xx}(0) + \mathbf{r}_{xx}^H(1) \mathbf{a}. \end{aligned} \quad (5.166)$$

**Autoregressive Processes and the Yule–Walker Equations.** We consider an autoregressive process of order  $p$  (AR( $p$ ) process). As outlined in Section 5.3, such a process is generated by exciting a stable recursive filter with a stationary white noise process  $w(n)$ . The system function of the recursive system is supposed to be<sup>2</sup>

$$U(z) = \frac{1}{1 + \sum_{i=1}^p a(i) z^{-i}}, \quad a(p) \neq 0. \quad (5.167)$$

The input-output relation of the recursive system may be expressed via the difference equation

$$x(n) = w(n) - \sum_{i=1}^p a(i) x(n-i). \quad (5.168)$$

---

<sup>2</sup>In order to keep in line with the notation used in the literature, the coefficients  $\rho(i)$ ,  $i = 1, \dots, p$  introduced in (5.34) are replaced by the coefficients  $-a(i)$ ,  $i = 1, \dots, p$ .

For the autocorrelation sequence of the process  $x(n)$  we thus derive

$$\begin{aligned} r_{xx}(m) &= E \{x^*(n)x(n+m)\} \\ &= r_{xw}(m) - \sum_{i=1}^p a(i) r_{xx}(m-i). \end{aligned} \quad (5.169)$$

The cross correlation sequence  $r_{xw}(m)$  is

$$\begin{aligned} r_{xw}(m) &= E \{x^*(n)w(n+m)\} \\ &= \sum_{i=1}^{\infty} u^*(i) \underbrace{r_{ww}(i+m)}_{\sigma_w^2 \delta(i+m)} \\ &= \sigma_w^2 u^*(-m), \end{aligned} \quad (5.170)$$

where  $u(n)$  is the impulse response of the recursive filter. Since  $u(n)$  is causal ( $u(n) = 0$  for  $n < 0$ ), we derive

$$r_{xw}(m) = \begin{cases} 0, & m < 0, \\ \sigma_w^2 u^*(-m), & m \geq 0. \end{cases} \quad (5.171)$$

By combining (5.169) and (5.171) we finally get

$$r_{xx}(m) = \begin{cases} -\sum_{i=1}^p a(i) r_{xx}(m-i), & m > 0, \\ \sigma_w^2 - \sum_{i=1}^p a(i) r_{xx}(m-i), & m = 0, \\ r_{xx}^*(-m), & m < 0. \end{cases} \quad (5.172)$$

The equations (5.172) are known as the *Yule-Walker equations*. In matrix form they are

$$\begin{bmatrix} r_{xx}(0) & r_{xx}(-1) & r_{xx}(-2) & \cdots & r_{xx}(-p) \\ r_{xx}(1) & r_{xx}(0) & r_{xx}(-1) & \cdots & r_{xx}(1-p) \\ \vdots & \vdots & \vdots & & \vdots \\ r_{xx}(p) & r_{xx}(p-1) & r_{xx}(p-1) & \cdots & r_{xx}(0) \end{bmatrix} \begin{bmatrix} 1 \\ a(1) \\ \vdots \\ a(p) \end{bmatrix} = \begin{bmatrix} \sigma_w^2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (5.173)$$

As can be inferred from (5.173), we obtain the coefficients  $a(i)$ ,  $i = 1, \dots, p$  by solving (5.163). By observing the power of the prediction error we can also determine the power of the input process. From (5.166) and (5.172) we have

$$\begin{aligned} \sigma_w^2 &= \sigma_{e_{\min}}^2 \\ &= r_{xx}(0) + \mathbf{r}_{xx}^H(1) \mathbf{a}. \end{aligned} \quad (5.174)$$

Thus, all parameters of an autoregressive process can be exactly determined from the parameters of a one-step linear predictor.

**Prediction Error Filter.** The output signal of the so-called *prediction error filter* is the signal  $e(n)$  in Figure 5.4 with the coefficients  $a(n)$  according to (5.163). Introducing the coefficient  $a(0) = 1$ ,  $e(n)$  is given by

$$e(n) = \sum_{i=0}^p a(i) x(n-i), \quad a(0) = 1. \quad (5.175)$$

The system function of the prediction error filter is

$$A(z) = 1 + \sum_{i=1}^p a(i)z^{-i} = \sum_{i=0}^p a(i)z^{-i}, \quad a(0) = 1. \quad (5.176)$$

In the special case that  $x(n)$  is an autoregressive process, the prediction error filter  $A(z)$  is the inverse system to the recursive filter  $U(z) \longleftrightarrow u(n)$ . This also means that the output signal of the prediction error filter is a white noise process. Hence, the prediction error filter performs a whitening transform and thus constitutes an alternative to the methods considered in Section 5.4. If  $x(n)$  is not truly autoregressive, the whitening transform is carried out at least approximately.

**Minimum Phase Property of the Prediction Error Filter.** Our investigation of autoregressive processes showed that the prediction error filter  $A(z)$  is inverse to the recursive filter  $U(z)$ . Since a stable filter does not have poles outside the unit circle of the  $z$ -plane, the corresponding prediction error filter cannot have zeros outside the unit circle. Even if  $x(n)$  is not an autoregressive process, we obtain a minimum phase prediction error filter, because the calculation of  $A(z)$  only takes into account the second-order statistics, which do not contain any phase information, cf. (1.105).

### 5.6.3 Filter Design on the Basis of Finite Data Ensembles

In the previous sections we assumed stationary processes and considered the correlation sequences to be known. In practice, however, linear predictors must be designed on the basis of a finite number of observations.

In order to determine the predictor filter  $a(n)$  from measured data  $x(1), x(2), \dots, x(N)$ , we now describe the prediction error

$$e(n) = x(n) + \sum_{i=1}^p a(i)x(n-i)$$

via the following matrix equation:

$$\mathbf{e} = \mathbf{X} \mathbf{a} + \mathbf{x}, \quad (5.177)$$

where  $\mathbf{a}$  contains the predictor coefficients, and  $\mathbf{X}$  and  $\mathbf{x}$  contain the input data. The term  $\mathbf{X} \mathbf{a}$  describes the convolution of the data with the impulse response  $a(n)$ .

The criterion

$$\|\mathbf{e}\| = \|\mathbf{X} \mathbf{a} + \mathbf{x}\| \stackrel{!}{=} \min \quad (5.178)$$

leads to the following normal equation:

$$\mathbf{X}^H \mathbf{X} \mathbf{a} = -\mathbf{X}^H \mathbf{x}. \quad (5.179)$$

Here, the properties of the predictor are dependent on the definition of  $\mathbf{X}$  and  $\mathbf{x}$ . In the following, two relevant methods will be discussed.

**Autocorrelation Method.** The *autocorrelation method* is based on the following estimation of the autocorrelation sequence:

$$\hat{r}_{xx}^{(AC)}(m) = \frac{1}{N} \sum_{n=1}^{N-|m|} x^*(n) x(n+m). \quad (5.180)$$

As can be seen,  $\hat{r}_{xx}^{(AC)}(m)$  is a biased estimate of the true autocorrelation sequence  $r_{xx}(m)$ , which means that  $E\{\hat{r}_{xx}^{(AC)}(m)\} \neq r_{xx}(m)$ . Thus, the autocorrelation method yields a biased estimate of the parameters of an autoregressive process. However, the correlation matrix  $\hat{\mathbf{R}}_{xx}^{(AC)}$  built from  $\hat{r}_{xx}^{(AC)}(m)$  has a Toeplitz structure, which enables us to efficiently solve the equation

$$\hat{\mathbf{R}}_{xx}^{(AC)} \hat{\mathbf{a}} = -\hat{r}_{xx}^{(AC)}(1) \quad (5.181)$$

by means of the *Levinson–Durbin recursion* [89, 47] or the *Schur algorithm* [130]. Textbooks that cover this topic are, for instance, [84, 99, 117].

The autocorrelation method can also be viewed as the solution to the problem (5.178) with

$$\mathbf{x} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ x(N) \\ \vdots \\ x(p+1) \\ \vdots \\ x(2) \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} & & & x(N) \\ & & \dots & \vdots \\ & x(N) & \dots & x(N-p+1) \\ x(N-1) & \dots & x(N-p) & \\ \vdots & & \vdots & \\ x(p) & \dots & x(1) & \\ \vdots & \dots & & \\ x(1) & & & \end{bmatrix} \quad (5.182)$$

and

$$\mathbf{a} = [a(1), a(2), \dots, a(p)]^T. \quad (5.183)$$

We have

$$\hat{\mathbf{R}}_{xx}^{(AC)} = \mathbf{X}^H \mathbf{X} \quad (5.184)$$

and

$$\hat{\mathbf{r}}_{xx}^{(AC)}(1) = \mathbf{X}^H \mathbf{x}. \quad (5.185)$$

**Covariance Method.** The *covariance method* takes into account the prediction errors in steady state only and yields an unbiased estimate of the autocorrelation matrix. In this case  $\mathbf{X}$  and  $\mathbf{x}$  are defined as

$$\mathbf{X} = \begin{bmatrix} x(N-1) & \dots & x(N-p) \\ \vdots & & \vdots \\ x(p+1) & \dots & x(2) \\ x(p) & \dots & x(1) \end{bmatrix} \quad (5.186)$$

and

$$\mathbf{x} = \begin{bmatrix} x(N) \\ \vdots \\ x(p+2) \\ x(p+1) \end{bmatrix}. \quad (5.187)$$

The equation to be solved is

$$\hat{\mathbf{R}}_{xx}^{(CV)} \hat{\mathbf{a}} = -\hat{\mathbf{r}}_{xx}^{(CV)}(1), \quad (5.188)$$

where

$$\hat{\mathbf{R}}_{xx}^{(CV)} = \mathbf{X}^H \mathbf{X}, \quad (5.189)$$

$$\hat{\mathbf{r}}_{xx}^{(CV)}(1) = \mathbf{X}^H \mathbf{x}. \quad (5.190)$$

Note that  $\hat{\mathbf{R}}_{xx}^{(CV)}$  is not a Toeplitz matrix, so that solving (5.188) is much more complex than solving (5.181) via the Levinson–Durbin recursion. However, the covariance method has the advantage of being unbiased; we have

$$E \left\{ \hat{\mathbf{R}}_{xx}^{(CV)} \right\} = \mathbf{R}_{xx}. \quad (5.191)$$

## 5.7 Estimation of Autocorrelation Sequences and Power Spectral Densities

### 5.7.1 Estimation of Autocorrelation Sequences

In the following, we will discuss methods for estimating the autocorrelation sequence of random processes from given sample values  $x(n)$ ,  $n = 0, \dots, N-1$ . We start the discussion with the estimate

$$\hat{r}_{xx}^b(m) = \frac{1}{N} \sum_{n=0}^{N-|m|-1} x^*(n) x(n+m), \quad (5.192)$$

which is the same as the estimate  $\hat{r}_{xx}^{(AC)}(m)$ , used in the autocorrelation method explained in the last section. As can easily be verified, the estimate  $\hat{r}_{xx}^b(m)$  is biased with mean

$$E\{\hat{r}_{xx}^b(m)\} = \frac{N-|m|}{N} r_{xx}(m). \quad (5.193)$$

However, since

$$\lim_{N \rightarrow \infty} E\{\hat{r}_{xx}^b(m)\} = r_{xx}(m), \quad (5.194)$$

the estimate is asymptotically unbiased. The triangular window  $\frac{N-|m|}{N}$  that occurs in (5.193) is known as the *Bartlett window*.

The variance of the estimate can be approximated as [77]

$$\text{var}[\hat{r}_{xx}^b(m)] \approx \frac{1}{N} \sum_{n=-\infty}^{\infty} |r_{xx}(n)|^2 + r_{xx}^*(n-m) r_{xx}(n+m). \quad (5.195)$$

Thus, as  $N \rightarrow \infty$ , the variance tends to zero:

$$\lim_{N \rightarrow \infty} \text{var}[\hat{r}_{xx}^b(m)] \rightarrow 0. \quad (5.196)$$

Such an estimate is said to be *consistent*. However, although consistency is given, we cannot expect good estimates for large  $m$  as long as  $N$  is finite, because the bias increases as  $|m| \rightarrow N$ .

**Unbiased Estimate.** An unbiased estimate of the autocorrelation sequence is given by

$$\hat{r}_{xx}^u(m) = \frac{1}{N-|m|} \sum_{n=0}^{N-|m|-1} x^*(n) x(n+m). \quad (5.197)$$

We have

$$E\{\hat{r}_{xx}^u(m)\} = r_{xx}(m), \quad |m| \leq N-1. \quad (5.198)$$

The variance of the estimate can be approximated as [77]

$$\text{var}[\hat{r}_{xx}^b(m)] \approx \frac{N}{(N - |m|)^2} \sum_{n=-\infty}^{\infty} |r_{xx}(n)|^2 + r_{xx}^*(n - m) r_{xx}(n + m). \quad (5.199)$$

As  $N \rightarrow \infty$ , this gives

$$\lim_{N \rightarrow \infty} \text{var}[\hat{r}_{xx}^u(m)] \rightarrow 0, \quad (5.200)$$

which means that  $\hat{r}_{xx}^u(m)$  is a consistent estimate. However, problems arise for large  $m$  as long as  $N$  is finite, because the variance increases for  $|m| \rightarrow N$ .

## 5.7.2 Non-Parametric Estimation of Power Spectral Densities

In many real-world problems one is interested in knowledge about the power spectral density of the data to be processed. Typically, only a finite set of observations  $x(n)$  with  $n = 0, 1, \dots, N-1$  is available. Since the power spectral density is the Fourier transform of the autocorrelation sequence, and since we have methods for the estimation of the autocorrelation sequence, it is a logical consequence to look at the Fourier transforms of these estimates. We start with  $\hat{r}_{xx}^b(m)$ . The Fourier transform of  $\hat{r}_{xx}^b(m)$  will be denoted as

$$P_{xx}(e^{j\omega}) = \sum_{m=-(N-1)}^{N-1} \hat{r}_{xx}^b(m) e^{-j\omega m}. \quad (5.201)$$

We know that  $\hat{r}_{xx}^b(m)$  is a biased estimate of the true autocorrelation sequence  $r_{xx}(m)$ , so that we can conclude that the spectrum  $P_{xx}(e^{j\omega})$  is a biased estimate of the true power spectral density  $S_{xx}(e^{j\omega})$ . In order to be explicit, let us recall that

$$E\{\hat{r}_{xx}^b(m)\} = w_B(m) r_{xx}(m), \quad (5.202)$$

with  $w_B(m)$  being the Bartlett window; i.e.

$$w_B(m) = \frac{N - |m|}{N}. \quad (5.203)$$

In the spectral domain, we have

$$\begin{aligned}
 E \{P_{xx}(e^{j\omega})\} &= \sum_{m=-(N-1)}^{N-1} E \{\hat{r}_{xx}^b(m)\} e^{-j\omega m} \\
 &= \sum_{m=-(N-1)}^{N-1} w_B(m) r_{xx}(m) e^{-j\omega m} \quad (5.204) \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{xx}(e^{j\nu}) W_B(e^{j(\omega-\nu)}) d\nu,
 \end{aligned}$$

where  $W_B(e^{j\omega})$  is the Fourier transform of  $w_B(m)$  given by

$$W_B(e^{j\omega}) = \frac{1}{N} \left( \frac{\sin(\omega N/2)}{\sin(\omega/2)} \right)^2. \quad (5.205)$$

Thus,  $E \{P_{xx}(e^{j\omega})\}$  is a smoothed version of the true power spectral density  $S_{xx}(e^{j\omega})$ , where smoothing is carried out with the Fourier transform of the Bartlett window.

A second way of computing  $P_{xx}(e^{j\omega})$  is to compute the Fourier transform of  $x(n)$  first and to derive  $P_{xx}(e^{j\omega})$  from  $X(e^{j\omega})$ . By inserting (5.192) into (5.201) and rearranging the expression obtained, we get

$$P_{xx}(e^{j\omega}) = \frac{1}{N} |X(e^{j\omega})|^2. \quad (5.206)$$

In the form (5.206)  $P_{xx}(e^{j\omega})$  is known as the *periodogram*.

Another way of deriving an estimate of the power spectral density is to consider the Fourier transform of the estimate  $\hat{r}_{xx}^u(m)$ . We use the notation  $Q_{xx}(e^{j\omega})$  for this type of estimate:

$$Q_{xx}(e^{j\omega}) = \sum_{m=-(N-1)}^{N-1} \hat{r}_{xx}^u(m) e^{-j\omega m}. \quad (5.207)$$

The expected value is

$$\begin{aligned}
 E \{Q_{xx}(e^{j\omega})\} &= \sum_{m=-(N-1)}^{N-1} E \{\hat{r}_{xx}^u(m)\} e^{-j\omega m} \\
 &= \sum_{m=-(N-1)}^{N-1} r_{xx}(m) e^{-j\omega m} \\
 &= \sum_{m=-\infty}^{\infty} w_R(m) r_{xx}(m) e^{-j\omega m} \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{xx}(e^{j\nu}) W_R(e^{j(\omega-\nu)}) d\nu,
 \end{aligned} \tag{5.208}$$

where  $w_R(m)$  is the rectangular window

$$w_R(m) = \begin{cases} 1, & \text{for } |m| \leq N-1 \\ 0, & \text{otherwise,} \end{cases} \tag{5.209}$$

and  $W_R(e^{j\omega})$  is its Fourier transform:

$$W_R(e^{j\omega}) = \frac{\sin(\omega[2N-1]/2)}{\sin(\omega/2)}. \tag{5.210}$$

This means that although  $\hat{r}_{xx}^u(m)$  is an unbiased estimate of  $r_{xx}(m)$ , the quantity  $Q_{xx}(e^{j\omega})$  is a biased estimate of  $S_{xx}(e^{j\omega})$ . The reason for this is the fact that only a finite number of taps of the autocorrelation sequence is used in the computation of  $Q_{xx}(e^{j\omega})$ . The mean  $E \{Q_{xx}(e^{j\omega})\}$  is a smoothed version of  $S_{xx}(e^{j\omega})$ , where smoothing is carried out with the Fourier transform of the rectangular window.

As  $N \rightarrow \infty$  both estimates  $\hat{r}_{xx}^b(m)$  and  $\hat{r}_{xx}^u(m)$  become unbiased. The same holds for  $P_{xx}(e^{j\omega})$  and  $Q_{xx}(e^{j\omega})$ , so that both estimates of the power spectral density are asymptotically unbiased. The behavior of the variance of the estimates is different. While the estimates of the autocorrelation sequences are consistent, those of the power spectral density are not. For example, for a Gaussian process  $x(n)$  with power spectral density  $S_{xx}(e^{j\omega})$ , the variance of the periodogram becomes

$$\text{var} [P_{xx}(e^{j\omega})] = \left[ 1 + \left( \frac{\sin(\omega N)}{N \sin \omega} \right)^2 \right] S_{xx}^2(e^{j\omega}), \tag{5.211}$$

which yields

$$\lim_{N \rightarrow \infty} \text{var} [P_{xx}(e^{j\omega})] = S_{xx}^2(e^{j\omega}). \tag{5.212}$$

Thus, the periodogram does not give a consistent estimate of  $S_{xx}(e^{j\omega})$ . The proof of (5.211) is straightforward and is omitted here.

**Use of the DFT or FFT for Computing the Periodogram.** Since the periodogram is computed from the Fourier transform of the finite data sequence, it can be efficiently evaluated at a discrete set of frequencies by using the FFT. Given a length- $N$  sequence  $x(n)$ , we may consider a length- $N$  DFT, resulting in

$$P_{xx}(e^{j\omega_k}) = \frac{1}{N} \left| \sum_{n=0}^{N-1} x(n) e^{-j2\pi k/N} \right|^2 \quad (5.213)$$

with  $\omega_k = 2\pi k/N$ . In many applications, the obtained number of samples of  $P_{xx}(e^{j\omega})$  may be insufficient in order to draw a clear picture of the periodogram. Moreover, the DFT length may be inconvenient for computation, because no powerful FFT algorithm is at hand for the given length. These problems can be solved by extending the sequence  $x(n)$  with zeros to an arbitrary length  $N' \geq N$ . This procedure is known as *zero padding*. We obtain

$$P_{xx}(e^{j\omega_k}) = \frac{1}{N'} \left| \sum_{n=0}^{N-1} x(n) e^{-j2\pi k/N'} \right|^2 \quad (5.214)$$

with  $\omega_k = 2\pi k/N'$ . The evaluation of (5.214) is typically carried out via the FFT.

**Bartlett Method.** Various methods have been proposed for achieving consistent estimates of the power spectral density. The *Bartlett method* does this by decomposing the sequence  $x(n)$  into disjoint segments of smaller length and taking the ensemble average of the spectrum estimates derived from the smaller segments. With

$$x^{(i)}(n) = x(n + iM), \quad i = 0, 1, \dots, K-1, \quad n = 0, 1, \dots, M-1, \quad (5.215)$$

we get the  $K$  periodograms

$$P_{xx}^{(i)}(e^{j\omega}) = \frac{1}{M} \left| \sum_{n=0}^{M-1} x^{(i)}(n) e^{-j\omega n} \right|^2, \quad i = 0, 1, \dots, K-1. \quad (5.216)$$

The Bartlett estimate then is

$$P_{xx}^B(e^{j\omega}) = \frac{1}{K} \sum_{i=0}^{K-1} P_{xx}^{(i)}(e^{j\omega}). \quad (5.217)$$

The expected value becomes

$$E \left\{ P_{xx}^B(e^{j\omega}) \right\} = E \left\{ P_{xx}^{(i)}(e^{j\omega}) \right\} = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{xx}(e^{j\nu}) W_{B_M}(e^{j(\omega-\nu)}) d\nu, \quad (5.218)$$

with  $W_{BM}(e^{j\omega})$  being the Fourier transform of the length- $M$  Bartlett window. Assuming a Gaussian process  $x(n)$ , the variance becomes

$$\text{var}[P_{xx}^B(e^{j\omega})] = \frac{1}{K} \text{var}[P_{xx}(e^{j\omega})] = \frac{1}{K} \left[ 1 + \left( \frac{\sin(\omega M)}{M \sin \omega} \right)^2 \right] S_{xx}^2(e^{j\omega}). \quad (5.219)$$

Thus, as  $N, M, K \rightarrow \infty$ , the variance tends to zero and the estimate is consistent. For finite  $N$ , the decomposition of  $x(n)$  into  $K$  sets results in a reduced variance, but the bias increases accordingly and the spectrum resolution decreases.

**Blackman-Tukey Method.** Blackman and Tukey proposed windowing the estimated autocorrelation sequence prior the Fourier transform [8]. The argument is that windowing allows us to reduce the influence of the unreliable estimates of the autocorrelation sequence for large  $m$ . Denoting the window and its Fourier transform as  $w(m)$  and  $W(e^{j\omega})$ , respectively, the estimate can be written as

$$P_{xx}^{BT}(e^{j\omega}) = \sum_{m=-(N-1)}^{N-1} w(m) \hat{r}_{xx}^b(m) e^{-j\omega m}. \quad (5.220)$$

In the frequency domain, this means that

$$P_{xx}^{BT}(e^{j\omega}) = \frac{1}{2\pi} \int_{-\pi}^{\pi} P_{xx}(e^{j\nu}) W(e^{j(\omega - \nu)}) d\nu. \quad (5.221)$$

The window  $w(n)$  should be chosen such that

$$W(e^{j\omega}) > 0 \quad \forall \omega \quad (5.222)$$

in order to ensure that  $P_{xx}^{BT}(e^{j\omega})$  is positive for all frequencies.

The expected value of  $P_{xx}^{BT}(e^{j\omega})$  is most easily expressed in the form

$$E \left\{ P_{xx}^{BT}(e^{j\omega}) \right\} = \sum_{m=-(N-1)}^{N-1} w(m) w_B(m) r_{xx}(m) e^{-j\omega m}. \quad (5.223)$$

Provided that  $w(m)$  is wide with respect to  $r_{xx}(m)$  and narrow with respect to  $w_B(m)$ , the expected value can be approximated as

$$E \left\{ P_{xx}^{BT}(e^{j\omega}) \right\} = w(0) S_{xx}(e^{j\omega}). \quad (5.224)$$

Thus, in order to achieve an asymptotically unbiased estimate, the window should satisfy

$$w(0) = 1. \quad (5.225)$$

For a symmetric window  $w(m) = w(-m)$  the variance can be estimated as [8]

$$\text{var} [P_{xx}^{BT}(e^{j\omega})] \approx \frac{1}{N} \left[ \sum_{m=-(N-1)}^{N-1} w^2(n) \right] S_{xx}^2(e^{j\omega}). \quad (5.226)$$

This approximation is based on the assumption that  $W(e^{j\omega})$  is wide with respect to  $W_B(e^{j\omega})$  and narrow with respect to the variations of  $S_{xx}(e^{j\omega})$ .

**Welch Method.** In the *Welch method* [162] the data is divided into overlapping blocks

$$x^{(i)}(n) = x(n + iD), \quad i = 0, 1, \dots, K - 1, \quad n = 0, 1, \dots, M - 1 \quad (5.227)$$

with  $D \leq M$ . For  $D = M$  we approach the decomposition in the Bartlett method. For  $D < M$  we have more segments than in the Bartlett method.

Each block is windowed prior to computation of the periodogram, resulting in  $K$  spectral estimates

$$V_{xx}^{(i)}(e^{j\omega}) = \frac{1}{\alpha M} \left| \sum_{n=0}^{M-1} x^{(i)}(n) w(n) e^{-j\omega n} \right|^2, \quad i = 0, 1, \dots, K - 1. \quad (5.228)$$

The factor  $\alpha$  is chosen as

$$\alpha = \frac{1}{M} \sum_{m=0}^{M-1} w^2(m) = \frac{1}{M} \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{ww}^E(e^{j\omega}) d\omega, \quad (5.229)$$

which means that the analysis is carried out with a window of normalized energy. Taking the average yields the final estimate

$$P_{xx}^W(e^{j\omega}) = \frac{1}{K} \sum_{i=0}^{K-1} V_{xx}^{(i)}(e^{j\omega}). \quad (5.230)$$

The expected value becomes

$$E \left\{ P_{xx}^W(e^{j\omega}) \right\} = E \left\{ V_{xx}^{(i)}(e^{j\omega}) \right\} \quad (5.231)$$

with

$$E \left\{ V_{xx}^{(i)}(e^{j\omega}) \right\} = \frac{1}{\alpha M} \sum_{n=0}^{M-1} \sum_{m=0}^{M-1} w(n) w(m) r_{xx}(n-m) e^{-j\omega(n-m)}. \quad (5.232)$$

In the spectral domain, this can be rewritten as

$$E \left\{ V_{xx}^{(i)}(e^{j\omega}) \right\} = \frac{1}{\alpha M} \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{xx}(e^{j\nu}) S_{ww}^E(e^{j(\omega-\nu)}) d\nu, \quad (5.233)$$

where

$$S_{ww}^E(e^{j\omega}) = |W(e^{j\omega})|^2. \quad (5.234)$$

With increasing  $N$  and  $M$ ,  $S_{ww}^E(e^{j(\omega-\nu)})$  becomes narrow with respect to  $S_{xx}(e^{j\nu})$  and the expected value tends to

$$E\{V_{xx}^{(i)}(e^{j\omega})\} = \frac{1}{\alpha M} S_{xx}(e^{j\omega}) \frac{1}{2\pi} \int_{-\pi}^{\pi} S_{ww}^E(e^{j\omega}) d\omega = S_{xx}(e^{j\omega}). \quad (5.235)$$

This shows that the Welch method is asymptotically unbiased.

For a Gaussian process, the variance of the estimate is

$$\text{var}[P_{xx}^W(e^{j\omega})] = \frac{1}{K^2} \sum_{i=0}^{K-1} \sum_{j=0}^{K-1} E\{V_{xx}^{(i)}(e^{j\omega}) V_{xx}^{(j)}(e^{j\omega})\} - [E\{V_{xx}^{(i)}(e^{j\omega})\}]^2. \quad (5.236)$$

If no overlap is considered ( $D = M$ ), the expression reduces to

$$\text{var}[P_{xx}^W(e^{j\omega})] = \frac{1}{K} \text{var}[V_{xx}^{(i)}(e^{j\omega})] \approx \frac{1}{K} S_{xx}^2(e^{j\omega}). \quad (5.237)$$

For  $k \rightarrow \infty$  the variance approaches zero, which shows that the Welch method is consistent.

Various windows with different properties are known for the purpose of spectral estimation. In the following, a brief overview is given.

*Hanning Window.*

$$w(n) = \begin{cases} 0.5 - 0.5 \cos\left(\frac{2\pi n}{N-1}\right), & n = 0, 1, \dots, N-1 \\ 0, & \text{otherwise.} \end{cases} \quad (5.238)$$

*Hamming Window.*

$$w(n) = \begin{cases} 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right), & n = 0, 1, \dots, N-1 \\ 0, & \text{otherwise.} \end{cases} \quad (5.239)$$

*Blackman Window.*

$$w(n) = \begin{cases} 0.42 - 0.5 \cos\left(\frac{2\pi n}{N-1}\right) + 0.08 \cos\left(\frac{4\pi n}{N-1}\right), & n = 0, 1, \dots, N-1 \\ 0, & \text{otherwise.} \end{cases} \quad (5.240)$$

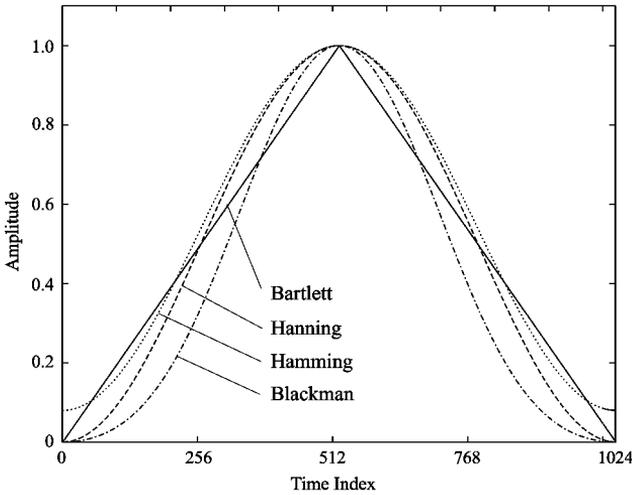


Figure 5.5. Window functions.

Figure 5.5 shows the windows, and Figure 5.6 shows their magnitude frequency responses. The spectrum of the Bartlett window is positive for all frequencies, which also means that the bias due to the Bartlett window is strictly positive. The spectra of the Hanning and Hamming window have relatively large negative side lobes, so that the estimated power spectral density may have a negative bias in the vicinity of large peaks in  $S_{xx}(e^{j\omega})$ . The Blackman window is a compromise between the Bartlett and the Hanning/Hamming approaches.

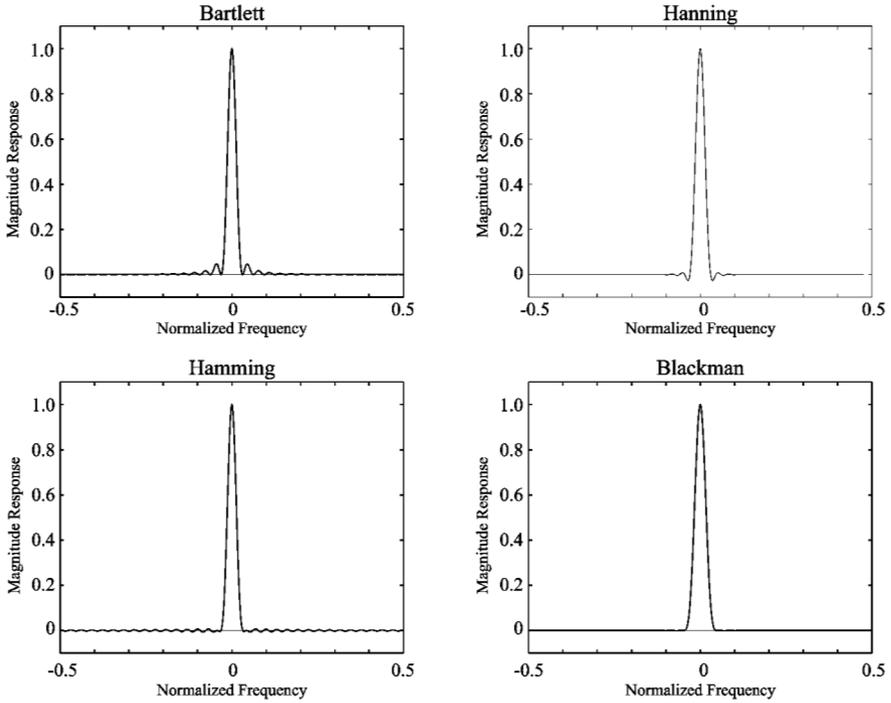
### 5.7.3 Parametric Methods in Spectral Estimation

Parametric methods in spectral estimation have been the subject of intensive research, and many different methods have been proposed. We will consider the simplest case only, which is related to the Yule–Walker equations. A comprehensive treatment of this subject would go far beyond the scope of this section.

Recall that in Section 5.6.2 we showed that the coefficients of a linear one-step predictor are identical to the parameters describing an autoregressive process. Hence the power spectral density may be estimated as

$$\hat{S}_{xx}(e^{j\omega}) = \frac{\hat{\sigma}_w^2}{\left| 1 + \sum_{n=1}^p \hat{a}(n)e^{-j\omega n} \right|^2}. \tag{5.241}$$

The coefficients  $\hat{a}(n)$  in (5.241) are the predictor coefficients determined from



**Figure 5.6.** Magnitude frequency responses of common window functions.

the observed data, and  $\hat{\sigma}_w^2$  is the power of the white input process estimated according to (5.174):

$$\hat{\sigma}_w^2 = \hat{r}_{xx}(0) + \hat{\mathbf{r}}_{xx}^H(1) \hat{\mathbf{a}}. \quad (5.242)$$

If we apply the autocorrelation method to the estimation of the predictor coefficients  $\hat{a}(n)$ , the estimated autocorrelation matrix has a Toeplitz structure, and the prediction filter is always minimum phase, just as when using the true correlation matrix  $\mathbf{R}_{xx}$ . For the covariance method this is not the case.

Finally, it shall be remarked that besides a forward prediction a backward prediction may also be carried out. By combining both predictors one can obtain an improved estimation of the power spectral density compared to (5.241). An example is the *Burg method* [19].