



EE210A: Adaptation and Learning

Professor Ali H. Sayed



MAIN RESULT FROM LAST LECTURE

Theorem 2.1 (Optimal estimation in the vector case) The least-mean-squares estimator of a (possibly complex-valued) vector x given another (possibly complex-valued) vector y is the conditional expectation of x given y , i.e., $\hat{x} = E(x|y)$. This estimator solves

$$\min_{\hat{x}} \text{Tr}(R_{\tilde{x}})$$

where $R_{\tilde{x}} = E\tilde{x}\tilde{x}^*$ and $\tilde{x} = x - \hat{x}$.

Lemma 2.2 (Cost function) The conditional expectation of x given y is optimal relative to either cost

$$\min_{\hat{x}} \text{Tr}(R_{\tilde{x}}) \quad \text{or} \quad \min_{\hat{x}} R_{\tilde{x}}$$

where $R_{\tilde{x}} = E\tilde{x}\tilde{x}^*$ and $\tilde{x} = x - \hat{x}$.

MAIN RESULT FROM LAST LECTURE

Lemma 2.1 (Circular Gaussian variables) If x and y are two circular and jointly Gaussian random variables with means $\{\bar{x}, \bar{y}\}$ and covariance matrices $\{R_x, R_y, R_{xy}\}$, then the least-mean-squares estimator of x given y is

$$\hat{x} = \bar{x} + R_{xy}R_y^{-1}(y - \bar{y})$$

and the resulting minimum cost is m.m.s.e. = $R_x - R_{xy}R_y^{-1}R_{yx}$.

$$R_x = E(x - \bar{x})(x - \bar{x})^*$$

$$R_y = E(y - \bar{y})(y - \bar{y})^*$$

$$R_{xy} = E(x - \bar{x})(y - \bar{y})^* = R_{yx}^*$$

LECTURE #03

LINEAR MEAN-SQUARE ERROR ESTIMATION

Sections in order: App. C, 3.1, 3.2, 3.3, 3.4, 4.2

COMPLEX GRADIENTS

C.1 CAUCHY-RIEMANN CONDITIONS



A. Cauchy G. Riemann
(1789-1857) (1826-1866)

We start with a scalar argument $z = x + jy$, where $j = \sqrt{-1}$. In this case, we can regard $g(z)$ as a function of the two real scalar variables, x and y , say

$$g(z) = u(x, y) + jv(x, y) \quad (\text{C.1})$$

with $u(\cdot, \cdot)$ denoting its real part and $v(\cdot, \cdot)$ denoting its imaginary part. Now, from complex function theory, the derivative of $g(z)$ at a point $z_o = x_o + jy_o$ is defined as

$$\frac{dg}{dz} \triangleq \lim_{\Delta z \rightarrow 0} \frac{g(x_o + \Delta x, y_o + \Delta y) - g(x_o, y_o)}{\Delta x + j\Delta y}$$

where $\Delta z = \Delta x + j\Delta y$. For $g(z)$ to be differentiable at z_o , in which case it is also said to be *analytic* at z_o , the above limit should exist regardless of the direction from which z approaches z_o . In particular, if we assume $\Delta y = 0$ and $\Delta x \rightarrow 0$, then the above definition gives

$$\frac{dg}{dz} = \frac{\partial u}{\partial x} + j \frac{\partial v}{\partial x} \quad (\text{C.2})$$

CAUCHY RIEMANN CONDITIONS

If, on the other hand, we assume that $\Delta x = 0$ and $\Delta y \rightarrow 0$ so that $\Delta z = j\Delta y$, then the definition gives

$$\frac{dg}{dz} = \frac{\partial v}{\partial y} - j \frac{\partial u}{\partial y} \quad (\text{C.3})$$

The expressions (C.2) and (C.3) should coincide. Therefore, by adding them we get

$$\frac{dg}{dz} = \frac{1}{2} \left(\frac{\partial u}{\partial x} + j \frac{\partial v}{\partial x} + \frac{\partial v}{\partial y} - j \frac{\partial u}{\partial y} \right)$$

or, more compactly,

$$\frac{dg}{dz} \triangleq \frac{1}{2} \left\{ \frac{\partial g}{\partial x} - j \frac{\partial g}{\partial y} \right\} \quad (\text{C.4})$$

Observe that the equality of expressions (C.2) and (C.3) implies that the real and imaginary parts of $g(\cdot)$ should satisfy the conditions

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} \quad \text{and} \quad \frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x}$$

SCALAR EXAMPLES

C.2 SCALAR ARGUMENTS

More generally, if g is a function of both z and z^* , we define its partial derivatives with respect to z and z^* as follows:

$$\frac{\partial g}{\partial z} = \frac{1}{2} \left\{ \frac{\partial g}{\partial x} - j \frac{\partial g}{\partial y} \right\}, \quad \frac{\partial g}{\partial z^*} = \frac{1}{2} \left\{ \frac{\partial g}{\partial x} + j \frac{\partial g}{\partial y} \right\} \quad (\text{C.5})$$

Note in particular, from the Cauchy-Riemann conditions, that if $g(z)$ is an analytic function of z then it must necessarily hold that $\partial g / \partial z^* = 0$.

Examples

We illustrate the definitions (C.4)–(C.5) considering several examples.

1. Let $g(z) = z = x + jy$. Then

$$\frac{\partial g}{\partial z} = (1 - j^2)/2 = 1, \quad \frac{\partial g}{\partial z^*} = (1 + j^2)/2 = 0$$

SCALAR EXAMPLES

2. Let $g(z) = z^2 = (x + jy)(x + jy) = (x^2 - y^2) + j2xy$. Then

$$\frac{\partial g}{\partial z} = 2(x + jy) = 2z, \quad \frac{\partial g}{\partial z^*} = 0$$

In examples 1. and 2., since g is a function of z alone, it holds that $\partial g / \partial z = dg/dz$.

3. Let $g(z) = |z|^2 = zz^* = (x + jy)(x - jy) = x^2 + y^2$. Then

$$\frac{\partial g}{\partial z} = (x - jy) = z^*, \quad \frac{\partial g}{\partial z^*} = (x + jy) = z$$

4. Let $g(z) = \lambda + \alpha z + \beta z^* + \gamma zz^*$, where $(\lambda, \alpha, \beta, \gamma)$ are complex constants. That is,

$$g(z) = [\lambda + \alpha x + \beta x + \gamma(x^2 + y^2)] + j[\alpha y - \beta y]$$

Then

$$\frac{\partial g}{\partial z} = \alpha + \gamma z^*, \quad \frac{\partial g}{\partial z^*} = \beta + \gamma z$$

VECTOR EXAMPLES

C.3 VECTOR ARGUMENTS

Now assume that z is a *column* vector, say

$$z = \text{col}\{z_1, z_2, \dots, z_n\}, \quad z_i = x_i + jy_i$$

The *complex gradient* of g with respect to z is denoted by $\nabla_z g(z)$, or simply $\nabla_z g$, and is defined as the *row* vector

$$\nabla_z g \triangleq \begin{bmatrix} \partial g / \partial z_1 & \partial g / \partial z_2 & \dots & \partial g / \partial z_n \end{bmatrix} \quad \left\{ \begin{array}{l} z \text{ is a column} \\ \nabla_z g \text{ is a row} \end{array} \right.$$

Likewise, the complex gradient of g with respect to z^* is defined as the *column* vector

$$\nabla_{z^*} g \triangleq \begin{bmatrix} \partial g / \partial z_1^* \\ \partial g / \partial z_2^* \\ \vdots \\ \partial g / \partial z_n^* \end{bmatrix} \quad \left\{ \begin{array}{l} z^* \text{ is a row} \\ \nabla_z g \text{ is a column} \end{array} \right.$$

VECTOR EXAMPLES

The reason why we choose to define $\nabla_z g$ as a row vector and $\nabla_{z^*} g$ as a column vector is because the subsequent differentiation results will be consistent with what we are used to from the standard differentiation of functions of real-valued arguments. Let us again consider a few examples:

1. Let $g(z) = \alpha^* z$, where $\{\alpha, z\}$ are column vectors. Then $\nabla_z g = \alpha^*$ and $\nabla_{z^*} g = 0$.
2. Let $g(z) = z^* \beta$, where $\{\beta, z\}$ are column vectors. Then $\nabla_z g = 0$ and $\nabla_{z^*} g = \beta$.
3. Let $g(z) = z^* z$, where z is a column vector. Then $\nabla_z g = z^*$ and $\nabla_{z^*} g = z$.
4. Let $g(z) = \lambda + \alpha^* z + z^* \beta + z^* \Gamma z$, where λ is a scalar, $\{\alpha, \beta\}$ are column vectors and Γ is a matrix. Then $\nabla_z g = \alpha^* + z^* \Gamma$ and $\nabla_{z^*} g = \beta + \Gamma z$.

OPTIMAL MSE ESTIMATOR

In Secs. 1.2 and 2.1 we studied the problem of determining the optimal function $h(\cdot)$ that minimizes the mean-square error of estimating a random variable x from another random variable y . Specifically, we solved

$$\min_{h(\cdot)} \mathbb{E} \tilde{x} \tilde{x}^* \quad (3.1)$$

over all functions $h(\cdot)$ of y . The optimal solution was found to be the conditional expectation of x given y , i.e.,

$$\hat{x} = \mathbb{E}(x|y)$$

Such conditional expectations are generally hard to evaluate in closed-form, except in some special cases.

AFFINE ESTIMATOR

Due to the difficulty in evaluating $E(x|y)$ in general, it is common practice to restrict the choice of $h(\cdot)$ to the subclass of *affine* functions of y , i.e., to functions of the form

$$h(y) = Ky + b \quad (3.2)$$

for some matrix K and for some vector b to be determined. For general vector-valued

Affine functions of the form (3.2) are easier to implement than most nonlinear functions. For example, when K is a row vector, the product Ky amounts to an inner product. Moreover, by deliberately restricting $h(\cdot)$ to be an affine function of y , the evaluation of the optimal $h(\cdot)$ is greatly simplified. In particular, it will be seen that all we need to know in order to determine the optimal parameters $\{K, b\}$ are the first and second-order moments of $\{x, y\}$, namely, $E x$, $E y$, $E xx^*$, $E yy^*$, and $E xy^*$. No other moments are needed. In contrast, evaluation of the optimal estimator, $\hat{x} = E(x|y)$, requires full knowledge of the conditional pdf $f_{x|y}(x|y)$.

3.1 MEAN-SQUARE ERROR CRITERION

Thus consider two zero-mean vector-valued random variables \mathbf{x} and \mathbf{y} and let

$$\bar{\mathbf{x}} \triangleq \mathbb{E}\mathbf{x} = \mathbf{0}, \quad \bar{\mathbf{y}} \triangleq \mathbb{E}\mathbf{y} = \mathbf{0}, \quad R_x \triangleq \mathbb{E}\mathbf{x}\mathbf{x}^*, \quad R_y \triangleq \mathbb{E}\mathbf{y}\mathbf{y}^*, \quad R_{xy} \triangleq \mathbb{E}\mathbf{x}\mathbf{y}^*$$

The dimensions of \mathbf{x} and \mathbf{y} need not be identical, say \mathbf{x} is $p \times 1$ and \mathbf{y} is $q \times 1$. In this case, $\{R_x, R_y, R_{xy}\}$ are $\{p \times p, q \times q, p \times q\}$ matrices, respectively.

We now seek an affine estimator for \mathbf{x} , namely, one of the form

$$\hat{\mathbf{x}} = K\mathbf{y} + b$$

for some constants $\{K, b\}$ to be determined, where K is $p \times q$ and b is $p \times 1$. The determination of $\{K, b\}$ is based on two considerations.

UNBIASED ESTIMATOR

The determination of $\{K, b\}$ is based on two considerations. First, the estimator should be unbiased, which means that it should satisfy

$$\mathbb{E} \hat{x} = 0$$

But since

$$\mathbb{E} \hat{x} = K \mathbb{E} y + b = 0 + b = b$$

we find that we must have $b = 0$. This means that the estimator that we are seeking is effectively a linear estimator, i.e., it is one of the form $\hat{x} = Ky$. Second, the coefficient matrix K should be chosen optimally so as to minimize the error covariance matrix (or its trace), as we now explain.

MINIMIZING MSE

Let $\{k_i^*\}$ denote the individual rows of K . Then the estimator for each entry of x , say $x(i)$, is given by the inner product k_i^*y ,

$$\underbrace{\begin{bmatrix} \hat{x}(0) \\ \hat{x}(1) \\ \vdots \\ \hat{x}(p-1) \end{bmatrix}}_{=\hat{x}} = \begin{bmatrix} k_0^*y \\ k_1^*y \\ \vdots \\ k_{p-1}^*y \end{bmatrix} = \underbrace{\begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{bmatrix}}_{=K} y$$

The optimal choices for the column vectors $\{k_i\}$ are determined by solving

$$\min_{k_i} \mathbb{E} |\tilde{x}(i)|^2 \quad \text{for each } i = 0, 1, \dots, p-1 \quad (3.3)$$

where

$$\tilde{x}(i) = x(i) - \hat{x}(i)$$

MINIMIZING MSE

The optimization problems (3.3) can be grouped together and stated equivalently as the problem of determining the matrix K by solving

$$\min_K \mathbb{E} \tilde{x}^* \tilde{x} \quad (3.4)$$

where $\tilde{x} = x - \hat{x}$. This is because the scalar quantity $\mathbb{E} \tilde{x}^* \tilde{x}$ in (3.4) is simply the sum of the individual error variances that appear in (3.3),

$$\mathbb{E} \tilde{x}^* \tilde{x} = \mathbb{E} |\tilde{x}(0)|^2 + \mathbb{E} |\tilde{x}(1)|^2 + \dots + \mathbb{E} |\tilde{x}(p-1)|^2 \quad (3.5)$$

and each term $\mathbb{E} |\tilde{x}(i)|^2$ depends on the corresponding k_i alone. In this way, minimizing $\mathbb{E} \tilde{x}^* \tilde{x}$ over K is equivalent to minimizing each term $\mathbb{E} |\tilde{x}(i)|^2$ over its k_i , so that problems (3.3) and (3.4) are equivalent.

MINIMIZING MSE

Therefore, continuing with (3.3), we expand the cost function to obtain a *quadratic* expression in the unknown column vector k_i ,

$$\begin{aligned}\mathsf{E}|\tilde{x}(i)|^2 &\stackrel{\Delta}{=} \mathsf{E}|x(i) - k_i^*y|^2 \\ &= \mathsf{E}[x(i) - k_i^*y][x(i) - k_i^*y]^* \\ &= \mathsf{E}|x(i)|^2 - [\mathsf{E}x(i)y^*]k_i - k_i^*[\mathsf{E}yx^*(i)] + k_i^*R_yk_i\end{aligned}$$

This is a *scalar-valued* cost function of a possibly *complex-valued vector* quantity, k_i . We denote it by

$$J(k_i) \stackrel{\Delta}{=} \mathsf{E}|x(i)|^2 - [\mathsf{E}x(i)y^*]k_i - k_i^*[\mathsf{E}yx^*(i)] + k_i^*R_yk_i \quad (3.6)$$

The quantity $\mathsf{E}|x(i)|^2$ that appears in the expression for $J(k_i)$ is the variance of $x(i)$ and is therefore equal to the i -th diagonal entry of R_x . We denote it by

$$\sigma_{x,i}^2 = \mathsf{E}|x(i)|^2$$

COST FUNCTION

Likewise, the quantity $\mathbb{E} \mathbf{x}(i)\mathbf{y}^*$ is the i -th row of the cross-covariance matrix R_{xy} . We denote it by $R_{xy,i} = \mathbb{E} \mathbf{x}(i)\mathbf{y}^*$. In this way, we can rewrite $J(k_i)$ as

$$J(k_i) \triangleq \sigma_{x,i}^2 - R_{xy,i}k_i - k_i^*R_{yx,i} + k_i^*R_yk_i \quad (3.7)$$

where $\{\sigma_{x,i}^2, R_{xy,i}, R_y\}$ are known quantities and k_i is the unknown column vector; $\sigma_{x,i}^2$ is a scalar, $R_{xy,i}$ is a row vector, and R_y is a nonnegative-definite matrix. Moreover, $R_{yx,i} = R_{xy,i}^*$. Our objective is to minimize $J(k_i)$ over k_i .

3.2 MINIMIZATION BY DIFFERENTIATION

$$J(k_i) \triangleq \sigma_{x,i}^2 - R_{xy,i}k_i - k_i^*R_{yx,i} + k_i^*R_yk_i \quad (3.7)$$

the complex gradient vector of $J(k_i)$ with respect to k_i is given by

$$\nabla_{k_i} J(k_i) = -R_{xy,i} + k_i^*R_y$$

Observe that this result is consistent with what we would expect from the standard rules of differentiation for functions of real variables, with the vectors k_i and k_i^* treated as different quantities. Also, if all data were real-valued, in which case

$$J(k_i) = \mathbb{E} x^2(i) - [\mathbb{E} x(i)\mathbf{y}^\top]k_i - k_i^\top[\mathbb{E} \mathbf{y}x(i)] + k_i^\top R_y k_i$$

then we would have obtained instead $\nabla_{k_i} J(k_i) = -R_{xy,i} + 2k_i^\top R_y$, with an additional factor of 2 — see App. C.

DIFFERENTIATION ARGUMENT

By setting the complex gradient equal to zero at the optimal choice $k_i = k_i^o$, we find that k_i^o should satisfy the linear equations

$$k_i^{o*} R_y = R_{xy,i}, \quad i = 0, 1, \dots, p-1 \quad (3.8)$$

If we collect the row vectors $\{k_i^{o*}\}$ from (3.8) into a matrix K_o we find that this desired solution matrix should satisfy

$$K_o R_y = R_{xy} \quad (3.9)$$

These equations are called the normal equations, for reasons explained later in Remark 4.1 in the next chapter.

COMPLETION OF SQUARES ARGUMENT

3.3 MINIMIZATION BY COMPLETION OF SQUARES

$$J(k_i) = \begin{bmatrix} 1 & k_i^* \end{bmatrix} \begin{bmatrix} \sigma_{x,i}^2 & -R_{xy,i} \\ -R_{yx,i} & R_y \end{bmatrix} \begin{bmatrix} 1 \\ k_i \end{bmatrix} \quad (3.10)$$

with a Hermitian center matrix and with the unknown vector k_i , and its conjugate transpose, multiplying from both sides.

Now given any Hermitian matrix of the form

$$M = \begin{bmatrix} A & B \\ B^* & C \end{bmatrix}$$

with $A = A^*$, $C = C^*$, and C invertible, it can be verified by direct calculation that M can be factored into a product of a block upper-triangular, block-diagonal, and block lower-triangular matrices as:

COMPLETION OF SQUARES

$$\begin{bmatrix} A & B \\ B^* & C \end{bmatrix} = \begin{bmatrix} I & BC^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & C \end{bmatrix} \begin{bmatrix} I & 0 \\ C^{-1}B^* & I \end{bmatrix} \quad (3.11)$$

where

$$\Sigma \triangleq A - BC^{-1}B^*$$

there is a generalization of (3.11) for block matrices M with possibly singular matrices C . Indeed, it is easy to verify, also by direct calculation, that we can alternatively factor any such M as

$$\begin{bmatrix} A & B \\ B^* & C \end{bmatrix} = \begin{bmatrix} I & D \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & C \end{bmatrix} \begin{bmatrix} I & 0 \\ D^* & I \end{bmatrix} \quad (3.12)$$

where D is defined as *any* solution to the linear system of equations

$$DC = B \quad (3.13)$$

and

$$\Sigma = A - BD^*$$

COMPLETION OF SQUARES

Applying (3.12) to the center matrix in (3.10) we can write

$$\begin{bmatrix} \sigma_{x,i}^2 & -R_{xy,i} \\ -R_{yx,i} & R_y \end{bmatrix} = \begin{bmatrix} 1 & -k_i^{o*} \\ 0 & I \end{bmatrix} \begin{bmatrix} \sigma_{x,i}^2 - R_{xy,i}k_i^o & 0 \\ 0 & R_y \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -k_i^o & I \end{bmatrix} \quad (3.14)$$

where k_i^o is any solution to the linear system of equations

$$k_i^{o*} R_y = R_{xy,i} \quad (3.15)$$

Substituting the factorization (3.14) into expression (3.10) for $J(k_i)$ and expanding the right-hand side, we find that $J(k_i)$ can be expressed in the equivalent form

$$J(k_i) = (\sigma_{x,i}^2 - R_{xy,i}k_i^o) + (k_i - k_i^o)^* R_y (k_i - k_i^o) \quad (3.16)$$

$$J(k_i^o) = \sigma_{x,i}^2 - R_{xy,i}k_i^o = \sigma_{x,i}^2 - k_i^{o*} R_y k_i^o = \sigma_{x,i}^2 - k_i^{o*} R_{yx,i} \quad (3.17)$$

$$\mathbb{E} \tilde{x}^* \tilde{x} = \text{Tr}(R_x - K_o R_y K_o^*) \quad (3.18)$$

ERROR COVARIANCE MATRIX

3.4 MINIMIZATION OF THE ERROR COVARIANCE MATRIX

The same completion-of-squares argument can be applied directly to the solution of (3.4) rather than the solution of the individual problems (3.3). To see this, let $J(K)$ denote the error-covariance matrix,

$$J(K) \triangleq \mathbb{E} \tilde{x} \tilde{x}^* = \mathbb{E} (x - Ky)(x - Ky)^*$$

so that, as in (3.10),

$$J(K) = [I \quad K] \begin{bmatrix} R_x & -R_{xy} \\ -R_{yx} & R_y \end{bmatrix} \begin{bmatrix} I \\ K^* \end{bmatrix}$$

ERROR COVARIANCE MATRIX

Following the same arguments as in the previous section, we can factor the center matrix as

$$\begin{bmatrix} R_x & -R_{xy} \\ -R_{yx} & R_y \end{bmatrix} = \begin{bmatrix} I & -K_o \\ 0 & I \end{bmatrix} \begin{bmatrix} R_x - R_{xy}K_o^* & 0 \\ 0 & R_y \end{bmatrix} \begin{bmatrix} I & 0 \\ -K_o^* & I \end{bmatrix}$$

where K_o is any solution to $K_o R_y = R_{xy}$. It then follows that

$$J(K) = (R_x - R_{xy}K_o^*) + (K - K_o)R_y(K - K_o)^* \quad (3.19)$$

where the last term is nonnegative definite for any K since $R_y \geq 0$. Now the criterion in (3.4) is related to $J(K)$ via $E \tilde{x}^* \tilde{x} = \text{Tr}[J(K)]$ so that, from (3.19),

$$E \tilde{x}^* \tilde{x} \geq \text{Tr}(R_x - R_{xy}K_o^*) \quad \text{for any } K \quad (3.20)$$

This is because the trace of the nonnegative-definite matrix $(K - K_o)R_y(K - K_o)^*$ is nonnegative.

ERROR COVARIANCE MATRIX

This derivation reveals another aspect of the solution K_o . It not only minimizes the trace of the error covariance matrix, as in (3.4), but it also minimizes the error-covariance matrix itself since we also get $J(K) \geq J(K_o)$ for any K . We can therefore interpret K_o as also the solution to the problem

$$\min_K \mathbb{E} \tilde{x}\tilde{x}^* \quad (3.21)$$

which is in terms of the error-covariance matrix, rather than its trace. The resulting minimum value is

$$J(K_o) = R_x - R_{xy}K_o^* \quad (3.22)$$

The optimization problem (3.21) is interesting for two reasons. First, the cost function $J(K) = \mathbb{E} \tilde{x}\tilde{x}^*$ is matrix-valued. That is, it assumes *matrix values* for each choice of K . Second, the unknown argument, K , is a *matrix* itself. In this way, problem (3.21) involves minimizing a matrix-valued cost function over a matrix-valued argument.

OPTIMAL LINEAR ESTIMATOR

Theorem 3.1 (Optimal linear estimator) Given zero-mean random variables x and y , the linear least-mean-squares estimator (l.l.m.s.e.) of x given y is

$$\hat{x} = K_o y$$

where K_o is any solution to the linear system of equations $K_o R_y = R_{xy}$. This estimator minimizes the following two error measures:

$$\min_K \mathbb{E} \tilde{x}^* \tilde{x} \quad \text{and} \quad \min_K \mathbb{E} \tilde{x} \tilde{x}^*$$

The scalar cost on the left is the trace of the matrix cost on the right. The resulting minimum mean-square errors, as defined by (3.4) and (3.22), are given by

$$\min_K \mathbb{E} \tilde{x}^* \tilde{x} = \text{Tr}(R_x - K_o R_y K_o^*)$$

$$\min_K \mathbb{E} \tilde{x} \tilde{x}^* = R_x - K_o R_y K_o^*$$

ORTHOGONALITY CONDITION

4.2 ORTHOGONALITY CONDITION

The linear least-mean-squares estimator admits an important geometric interpretation in the form of an orthogonality condition. This can be seen by rewriting the normal equations (3.9) as $K_o \mathbf{E} \mathbf{y} \mathbf{y}^* = \mathbf{E} \mathbf{x} \mathbf{y}^*$ or, equivalently,

$$\mathbf{E} (\mathbf{x} - K_o \mathbf{y}) \mathbf{y}^* = 0 \quad (4.3)$$

The difference $\mathbf{x} - K_o \mathbf{y}$ is the estimation error, $\tilde{\mathbf{x}}$. Therefore, equality (4.3) states that the error is orthogonal to (or uncorrelated with) the observation vector \mathbf{y} , namely,

$$\mathbf{E} \tilde{\mathbf{x}} \mathbf{y}^* = 0$$

which we also write as (see Fig. 4.4):

$$\tilde{\mathbf{x}} \perp \mathbf{y} \quad (4.4)$$

ORTHOGONALITY CONDITION

We thus conclude that for linear least-mean-squares estimation, the estimation error is orthogonal to the data and, in fact, to any *linear* transformation of the data, say Ay for any matrix A . This fact means that no further linear transformation of y can extract additional information about x in order to further reduce the error covariance matrix. Moreover, since the estimator \hat{x} is itself a linear function of y , we obtain, as a special case, that

$$\tilde{x} \perp \hat{x} \quad (4.5)$$

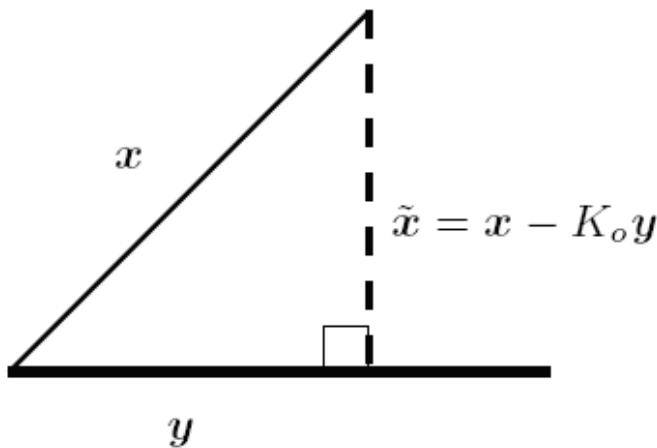


FIGURE 4.4 The orthogonality condition for linear estimation: $\tilde{x} \perp y$.

DEFINING PROPERTY

Theorem 4.1 (Orthogonality principle) Given two zero-mean random variables x and y , a linear estimator $\hat{x} = K_o y$ is optimal in the least-meansquares sense (3.3) if, and only if, it satisfies $x - \hat{x} \perp y$, i.e., $E(x - \hat{x})y^* = 0$.

Proof: One direction was argued prior to the statement of the theorem. Specifically, if \hat{x} is the optimal linear estimator, then we know from (4.4) that $\tilde{x} \perp y$. Conversely, assume \hat{x} is a linear estimator for x that satisfies $x - \hat{x} \perp y$, and let $\hat{x} = Ky$ for some K . It then follows from $x - \hat{x} \perp y$ that $E(x - Ky)y^* = 0$, so that K satisfies the normal equations $KR_y = R_{xy}$ and, hence, from Thm. 3.1 we conclude that \hat{x} should be the optimal linear estimator.



EXAMPLE

Example 4.5 (Signal with exponential auto-correlation)

Consider a scalar zero-mean stationary random process $\{z(t)\}$ with auto-correlation function:

$$R_z(\tau) \triangleq \mathbb{E} z(t)z^*(t - \tau) = e^{-\alpha|\tau|} \quad (4.6)$$

That is, the samples of $z(t)$ become less correlated as the time gap between them increases. A so-called *random telegraph* signal has this property — see Prob. II.23. It is claimed that the linear least-mean squares estimator of $z(T_3)$ given $z(T_1)$ and $z(T_2)$ (assuming $T_1 < T_2 < T_3$) is

$$\hat{z}(T_3) = e^{-\alpha(T_3-T_2)} z(T_2)$$

That is, the estimator of a future value depends only on the most recent observation, $z(T_2)$. We can verify the validity of this claim by checking whether the above estimator satisfies the orthogonality condition.

EXAMPLE

So define the observation vector $\mathbf{y} = \text{col}\{z(T_1), z(T_2)\}$ and let $x = z(T_3)$. We can now evaluate the cross-correlation vector

$$\mathbb{E}(x - \hat{x})\mathbf{y}^* = \mathbb{E}[z(T_3) - e^{-\alpha(T_3-T_2)}z(T_2)]\mathbf{y}^*$$

If the answer is zero then the orthogonality condition is satisfied and the estimator is optimal in the linear least-mean-squares sense, as claimed. Otherwise, the estimator is not optimal. Using the given auto-correlation function (4.6), it is easy to verify that

$$\begin{aligned}\mathbb{E}[z(T_3) - e^{-\alpha(T_3-T_2)}z(T_2)]\mathbf{y}^* &= \begin{bmatrix} e^{-\alpha(T_3-T_1)} - e^{-\alpha(T_3-T_1)} & e^{-\alpha(T_3-T_2)} - e^{-\alpha(T_3-T_2)} \end{bmatrix} \\ &= \begin{bmatrix} 0 & 0 \end{bmatrix}\end{aligned}$$

so that the estimator is optimal.



1st DESIGN EXAMPLE

Example 4.1 (Noisy measurements of a binary signal)

We reconsider Ex. 1.2 of a BPSK signal x that assumes the values ± 1 with probability $1/2$. The measurement y is $y = x + v$, where x and the disturbance v are independent of each other, with v being zero-mean Gaussian of unit variance.

Both x and y have zero means so that, according to Thm. 3.1, the optimal linear estimator of x is $\hat{x} = k_o y$, where the (now scalar) coefficient k_o is obtained from solving $k_o \sigma_y^2 = \sigma_{xy}$. We therefore need to determine the quantities $\{\sigma_y^2, \sigma_{xy}\}$. Now since $\{x, v\}$ are independent we have

$$\sigma_y^2 = \sigma_x^2 + \sigma_v^2 = 1 + 1 = 2$$

Moreover,

$$\sigma_{xy} = \mathbb{E} xy = \mathbb{E} x(x + v) = \mathbb{E} x^2 + 0 = \mathbb{E} x^2 = 1$$

so that $k_o = 1/2$, and the optimal linear estimator is $\hat{x} = y/2$. That is, we simply scale the received signal by $1/2$. In contrast, the optimal estimator was found in Ex. 1.2 to be given by the nonlinear transformation $\tanh(y)$. In addition, observe that the form of the linear estimator, $\hat{x} = y/2$, is valid regardless of whether the noise v is Gaussian or not (i.e., the Gaussian assumption on v is not needed to arrive at $k_o = 1/2$). The form of the optimal estimator, $\hat{x} = \tanh(y)$, on the other hand, is very much tied to the Gaussian assumption on v .

1st DESIGN EXAMPLE

Let us now reconsider Ex. 2.1, where we collect two noisy measurements $\mathbf{y}(0)$ and $\mathbf{y}(1)$ of x , say

$$\mathbf{y}(0) = \mathbf{x} + \mathbf{v}(0) \quad \text{and} \quad \mathbf{y}(1) = \mathbf{x} + \mathbf{v}(1)$$

where $\{\mathbf{v}(0), \mathbf{v}(1)\}$ are zero-mean unit-variance Gaussian random variables that are independent of each other and of \mathbf{x} . The value of \mathbf{x} is the same in both measurements (i.e., if it is $+1$ in the measurement $\mathbf{y}(0)$, it is also $+1$ in the measurement $\mathbf{y}(1)$, and similarly for -1) — recall Fig. 2.1. Introduce the column vector $\mathbf{y} = \text{col}\{\mathbf{y}(0), \mathbf{y}(1)\}$. Then, according to Thm. 3.1, the optimal linear estimator of \mathbf{x} given \mathbf{y} is $\hat{\mathbf{x}} = k_o^* \mathbf{y}$, where k_o^* is now 1×2 and is obtained from the solution of the normal equations $k_o^* R_y = R_{xy}$. To determine $\{R_y, R_{xy}\}$ we proceed as follows. Since $\{\mathbf{x}, \mathbf{v}(0), \mathbf{v}(1)\}$ are independent we get

$$R_y = \mathbb{E} \mathbf{y} \mathbf{y}^* = \begin{bmatrix} \mathbb{E} |\mathbf{y}(0)|^2 & \mathbb{E} \mathbf{y}(0) \mathbf{y}^*(1) \\ \mathbb{E} \mathbf{y}(1) \mathbf{y}^*(0) & \mathbb{E} |\mathbf{y}(1)|^2 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

1ST DESIGN EXAMPLE

where we used the fact that

$$\mathbb{E} \mathbf{y}(0)\mathbf{y}^*(1) = \mathbb{E} (\mathbf{x} + \mathbf{v}(0))(\mathbf{x} + \mathbf{v}(1))^* = \mathbb{E} |\mathbf{x}|^2 = 1$$

Likewise,

$$R_{xy} = \mathbb{E} \mathbf{x}\mathbf{y}^* = \begin{bmatrix} \mathbb{E} \mathbf{x}\mathbf{y}^*(0) & \mathbb{E} \mathbf{x}\mathbf{y}^*(1) \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

so that

$$k_o^* = R_{xy}R_y^{-1} = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} = \frac{1}{3} \begin{bmatrix} 1 & 1 \end{bmatrix}$$

That is,

$$\hat{\mathbf{x}} = \frac{1}{3} [\mathbf{y}(0) + \mathbf{y}(1)]$$

Again, this expression for the linear least-mean-squares estimator of \mathbf{x} given $\{\mathbf{y}(0), \mathbf{y}(1)\}$ holds regardless of whether the noises $\{\mathbf{v}(0), \mathbf{v}(1)\}$ are Gaussian or not. Only the first and second moments of $\{\mathbf{v}(0), \mathbf{v}(1)\}$, namely, their means and variances, are needed to determine k_o^* . In the context of the two-antenna example of Fig. 2.1, the above result leads to the optimal *linear* receiver structure shown in Fig. 4.1.

1ST DESIGN EXAMPLE

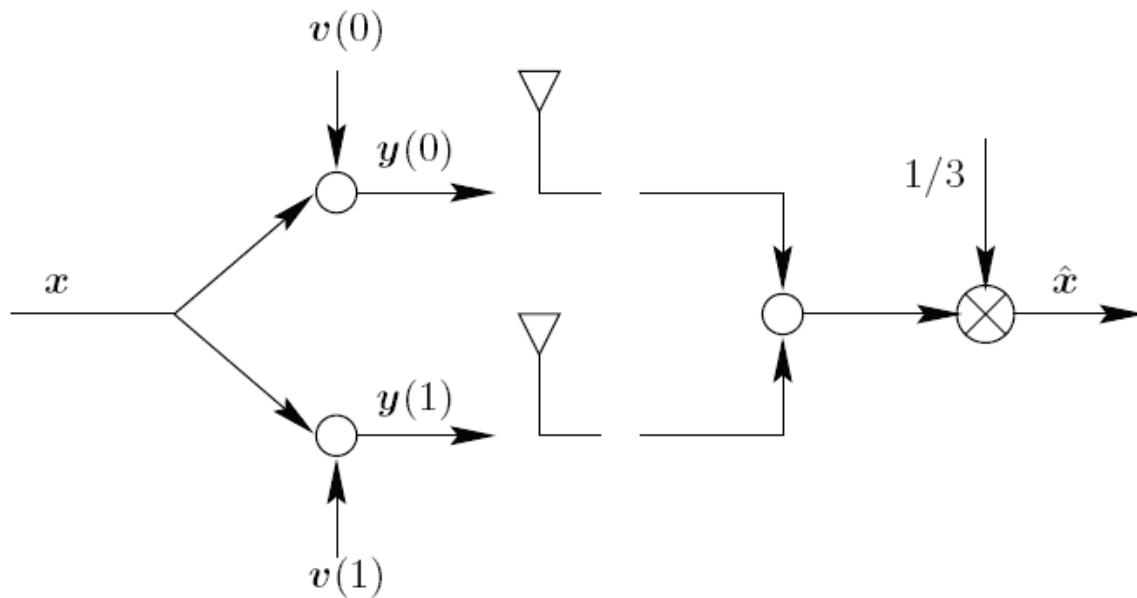


FIGURE 4.1 An optimal linear receiver for recovering a BPSK transmission from two measurements in the presence of additive unit-variance uncorrelated noises.

2ND DESIGN EXAMPLE

Example 4.2 (Multiple measurements of a binary signal)

Continuing with Ex. 4.1, let us examine what happens if we increase the number of available measurements from 2 to N , say

$$\mathbf{y}(i) = \mathbf{x} + \mathbf{v}(i), \quad i = 0, 1, \dots, N - 1$$

Introduce the observation vector $\mathbf{y} = \text{col}\{\mathbf{y}(0), \mathbf{y}(1), \dots, \mathbf{y}(N - 1)\}$. Then, say for $N = 5$,

$$R_{xy} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}, \quad R_y = \begin{bmatrix} 2 & 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & 2 \end{bmatrix}$$

2ND DESIGN EXAMPLE

so that

$$\hat{x} = R_{xy}R_y^{-1}\mathbf{y} = \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & 2 \end{bmatrix}^{-1} \right) \cdot \mathbf{y}$$

We need to evaluate R_y^{-1} . Due to the special structure of R_y , its inverse can be evaluated in closed form for any N , as we explain below. Later, in Sec. 5.5, when we reconsider this problem, we shall show how to evaluate \hat{x} via a more direct route.

Observe that, for any N , the matrix R_y can be expressed as $R_y = I + aa^T$, where a is the $N \times 1$ column vector $a = \text{col}\{1, 1, 1, \dots, 1\}$. In other words, R_y is a rank-one modification of the identity matrix. This is a useful observation since the inverse of every such matrix has a similar form (see Prob. II.1). Specifically,

2ND DESIGN EXAMPLE

$$(I + aa^T)^{-1} = I - \frac{aa^T}{1 + \|a\|^2} = I - \frac{1}{N+1}aa^T$$

where $\|a\|^2$ denotes the squared Euclidean norm of a , $\|a\|^2 = a^T a$. Using this result we find that

$$R_{xy}R_y^{-1} = a^T \left(I - \frac{aa^T}{N+1} \right) = a^T - \frac{N}{N+1}a^T = \frac{a^T}{N+1}$$

so that

$$\hat{x} = R_{xy}R_y^{-1}y = \frac{a^T}{N+1}y = \frac{1}{N+1} \sum_{k=0}^{N-1} y(k)$$

