

29 Linear Regression

The mean-square-error (MSE) problem of estimating a random variable \mathbf{x} from observations of another random variable \mathbf{y} seeks a mapping $c(\mathbf{y})$ that solves

$$\hat{\mathbf{x}} = \underset{\mathbf{x}=c(\mathbf{y})}{\operatorname{argmin}} \mathbb{E}(\mathbf{x} - c(\mathbf{y}))^2 \quad (29.1)$$

We showed in (27.18) that the optimal estimate is given by the conditional mean $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|\mathbf{y} = y)$. For example, for continuous random variables, the MSE estimate involves an integral computation of the form:

$$\hat{x} = \int_{x \in \mathcal{X}} x f_{\mathbf{x}|\mathbf{y}}(x|y) dx \quad (29.2)$$

over the domain of realizations $x \in \mathcal{X}$. Evaluation of this solution requires knowledge of the conditional distribution, $f_{\mathbf{x}|\mathbf{y}}(x|y)$. Even if $f_{\mathbf{x}|\mathbf{y}}(x|y)$ were available, computation of the integral expression is generally not possible in closed-form. In this chapter, we address this challenge by limiting $c(\mathbf{y})$ to the class of *affine* functions of \mathbf{y} . Assuming \mathbf{y} is M -dimensional, affine functions take the following form:

$$c(\mathbf{y}) = \mathbf{y}^\top \mathbf{w} - \theta \quad (29.3)$$

for some parameters $\mathbf{w} \in \mathbb{R}^M$ and $\theta \in \mathbb{R}$; the latter is called the *offset* parameter. The problem of determining the MSE estimator $\hat{\mathbf{x}}$ is then reduced to the problem of selecting optimal parameters (\mathbf{w}, θ) to minimize the same mean-square-error in (29.1). Despite its apparent narrowness, this class of estimators leads to solutions that are tractable mathematically and deliver laudable performance in a wide range of applications.

29.1 REGRESSION MODEL

Although we can treat the inference problem in greater generality than below, by considering directly the problem of estimating a random *vector* \mathbf{x} from another random *vector* \mathbf{y} , we will consider first the case of estimating a *scalar* x from a *vector* $\mathbf{y} \in \mathbb{R}^M$.

Let $\{\bar{x}, \bar{y}\}$ denote the first-order moments of the random variables $\mathbf{x} \in \mathbb{R}$ and

$\mathbf{y} \in \mathbb{R}^M$, i.e., their means:

$$\bar{x} = \mathbb{E} \mathbf{x}, \quad \bar{y} = \mathbb{E} \mathbf{y} \quad (29.4a)$$

and let $\{\sigma_x^2, R_y, r_{xy}\}$ denote their second-order moments, i.e., their (co)-variances and cross-covariance vector:

$$\sigma_x^2 = \mathbb{E} (\mathbf{x} - \bar{x})^2, \quad (\text{scalar}) \quad (29.4b)$$

$$R_y = \mathbb{E} (\mathbf{y} - \bar{y})(\mathbf{y} - \bar{y})^\top, \quad (M \times M) \quad (29.4c)$$

$$r_{xy} = \mathbb{E} (\mathbf{x} - \bar{x})(\mathbf{y} - \bar{y})^\top = r_{yx}^\top, \quad (1 \times M) \quad (29.4d)$$

The cross-covariance vector, r_{xy} , between \mathbf{x} and \mathbf{y} is a useful measure of the amount of information that one variable conveys about the other. We then pose the problem of determining the *linear* least-mean-square-error *estimator* (l.l.m.s.e., for short) of \mathbf{x} given \mathbf{y} , namely, an estimator for the form

$$\hat{\mathbf{x}}_{\text{LMSE}} = \mathbf{y}^\top w - \theta = w^\top \mathbf{y} - \theta, \quad (\text{estimator for } \mathbf{x}) \quad (29.5)$$

where (w, θ) are determined by solving

$$(w^\circ, \theta^\circ) = \underset{w, \theta}{\operatorname{argmin}} \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}}_{\text{LMSE}})^2 \quad (29.6)$$

The minus sign in front of the offset parameter θ in (29.5) is chosen for convenience, and the subscript LMSE refers to the “linear mean-square error estimator.” Since in this chapter we will be dealing almost exclusively with linear estimators under the mean-square-error criterion, we will refrain from including the subscript and will simply write $\hat{\mathbf{x}}$. It is customary to refer to model (29.5) as a *linear regression model* in the sense that the individual entries of \mathbf{y} are being combined linearly, or a linear model is being fitted to the entries of \mathbf{y} , in order to estimate \mathbf{x} .

THEOREM 29.1. (Linear estimators) *The solution (w°, θ°) to (29.6) satisfies the relations*

$$R_y w^\circ = r_{yx}, \quad \theta^\circ = \bar{y}^\top w^\circ - \bar{x} \quad (29.7)$$

so that, when R_y is invertible, the estimator and the resulting minimum mean-square error (m.m.s.e.) are given by

$$\hat{\mathbf{x}} - \bar{x} = r_{xy} R_y^{-1} (\mathbf{y} - \bar{y}) \quad (29.8a)$$

$$\text{m.m.s.e.} = \sigma_x^2 - r_{xy} R_y^{-1} r_{yx} \quad (29.8b)$$

Proof: We provide an algebraic proof. Let $\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$. We expand the mean-square error (m.s.e.) to get

$$\begin{aligned} \text{m.s.e.} &\triangleq \mathbb{E} \tilde{\mathbf{x}}^2 = \mathbb{E} (\mathbf{x} - \mathbf{y}^\top w + \theta)^2 \\ &= \mathbb{E} \mathbf{x}^2 - 2(\mathbb{E} \mathbf{x} \mathbf{y}^\top) w + 2\theta \bar{x} - 2\theta \bar{y}^\top w + w^\top (\mathbb{E} \mathbf{y} \mathbf{y}^\top) w + \theta^2 \end{aligned} \quad (29.9)$$

This m.s.e. is a quadratic function of w and θ . Differentiating with respect to w and θ and setting the derivatives to zero at the optimal solution gives:

$$\left. \frac{\partial \text{m.s.e.}}{\partial \theta} \right|_{\theta=\theta^o, w=w^o} = 2\bar{x} - 2\bar{y}^\top w^o + 2\theta^o = 0 \quad (29.10)$$

$$\left. \nabla_w \text{m.s.e.} \right|_{\theta=\theta^o, w=w^o} = -2\mathbb{E} \mathbf{xy}^\top - 2\theta^o \bar{y}^\top + 2(w^o)^\top (\mathbb{E} \mathbf{yy}^\top) = 0 \quad (29.11)$$

Solving for θ^o and w^o we find

$$\theta^o = \bar{y}^\top w^o - \bar{x} \quad \text{and} \quad (\mathbb{E} \mathbf{yy}^\top) w^o = \mathbb{E} \mathbf{yx} + \theta^o \bar{y} \quad (29.12)$$

Replacing θ^o in the second expression for w^o and grouping terms gives

$$(\mathbb{E} \mathbf{yy}^\top - \bar{y} \bar{y}^\top) w^o = \mathbb{E} \mathbf{yx} - \bar{y} \bar{x} \iff R_y w^o = r_{yx} \quad (29.13)$$

which leads to (29.7). Substituting the expressions for (w^o, θ^o) into (29.9) leads to (29.8b). Finally, the Hessian matrix of the m.s.e. relative to w and θ is given by

$$H \triangleq \begin{bmatrix} \frac{\partial^2 \text{m.s.e.}}{\partial \theta^2} & \frac{\partial}{\partial \theta} (\nabla_w \text{m.s.e.}) \\ \nabla_w^\top (\frac{\partial \text{m.s.e.}}{\partial \theta}) & \nabla_w^2 \text{m.s.e.} \end{bmatrix} = \begin{bmatrix} 2 & -2\bar{y}^\top \\ -2\bar{y} & 2\mathbb{E} \mathbf{yy}^\top \end{bmatrix} \quad (29.14)$$

This Hessian matrix is nonnegative-definite since its (1,1) entry is positive and the Schur complement relative to this entry is nonnegative-definite:

$$\text{Schur complement} = 2\mathbb{E} \mathbf{yy}^\top - 2\bar{y} \bar{y}^\top = 2R_y \geq 0 \quad (29.15)$$

Since the m.s.e. cost is quadratic in the parameters (w, θ) , we conclude that the solution (w^o, θ^o) corresponds to a global minimizer. ■

It is sufficient for our purposes to assume that $R_y > 0$. Observe from the statement of the theorem that the solution to the *linear* regression problem only requires knowledge of the first and second-order moments $\{\bar{x}, \bar{y}, \sigma_x^2, R_y, r_{xy}\}$; there is no need to know the full conditional pdf $f_{x|y}(x|y)$ as was the case with the optimal conditional mean estimator (27.18). Moreover, the observation, \mathbf{y} , does not appear in the m.m.s.e. expression. This means that we can assess beforehand, even before receiving the observation, the performance level that will be expected from the solution.

REMARK 29.1. (Linear model) Consider two random variables $\{\mathbf{x}, \mathbf{y}\}$ and assume they are related by a linear model of the form:

$$\mathbf{x} = \mathbf{y}^\top z^o + \mathbf{v} \quad (29.16)$$

for some unknown parameter $z^o \in \mathbb{R}^M$, and where \mathbf{v} has zero mean and is orthogonal to \mathbf{y} , i.e., $\mathbb{E} \mathbf{vy}^\top = 0$. Taking expectation of both sides gives $\bar{x} = \bar{y}^\top z^o$ so that

$$\mathbf{x} - \bar{x} = (\mathbf{y} - \bar{y})^\top z^o + \mathbf{v} \quad (29.17)$$

Multiplying both sides by $(\mathbf{y} - \bar{y})$ from the left and taking expectations again gives

$$r_{yx} = R_y z^o \quad (29.18)$$

This is the same equation satisfied by the solution w^o in (29.7). We therefore conclude that when the variables $\{\mathbf{x}, \mathbf{y}\}$ happen to be related by a linear model as in (29.16),

then the estimator, $\hat{\mathbf{x}}_{\text{LMSE}}$, is able to recover the *exact* model z° . Put in another way, when we solve a linear mean-square error problem, we are implicitly assuming that the variables $\{\mathbf{x}, \mathbf{y}\}$ satisfy a linear model of the above form. ■

Two properties

The linear least-mean-squares estimator satisfies two useful properties. First, the estimator is unbiased since by taking expectations of both sides of (29.8a) we get

$$\mathbb{E}(\hat{\mathbf{x}} - \bar{x}) = r_{xy} R_y^{-1} \underbrace{\mathbb{E}(\mathbf{y} - \bar{y})}_{=0} = 0 \quad (29.19)$$

so that,

$$\boxed{\mathbb{E} \hat{\mathbf{x}} = \bar{x}} \quad (\text{unbiased estimator}) \quad (29.20)$$

Moreover, using (29.8a) again, we find that the variance of the estimator $\hat{\mathbf{x}}$ is given by

$$\begin{aligned} \sigma_{\hat{x}}^2 &= \mathbb{E}(\hat{\mathbf{x}} - \bar{x})^2 \\ &\stackrel{(29.8a)}{=} r_{xy} R_y^{-1} \underbrace{\left(\mathbb{E}(\mathbf{y} - \bar{y})(\mathbf{y} - \bar{y})^\top \right)}_{=R_y} R_y^{-1} r_{yx} \\ &= r_{xy} R_y^{-1} r_{yx} \end{aligned} \quad (29.21)$$

so that expression (29.8b) for the m.m.s.e. can be equivalently rewritten as

$$\text{m.m.s.e.} = \mathbb{E} \tilde{\mathbf{x}}^2 = \sigma_x^2 - \sigma_{\hat{x}}^2 \quad (29.22)$$

Second, the linear least-mean-square error estimator satisfies an important orthogonality condition, namely, it is uncorrelated with the observation:

$$\boxed{\mathbb{E} \tilde{\mathbf{x}} \mathbf{y}^\top = 0} \quad (\text{orthogonality principle}) \quad (29.23)$$

Proof of (29.23): From expression (29.8a) we note that

$$\begin{aligned} \mathbb{E} \tilde{\mathbf{x}}(\mathbf{y} - \bar{y})^\top &= \mathbb{E}(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{y} - \bar{y})^\top \\ &\stackrel{(29.8a)}{=} \mathbb{E}(\mathbf{x} - \bar{x} - r_{xy} R_y^{-1}(\mathbf{y} - \bar{y}))(\mathbf{y} - \bar{y})^\top \\ &= \mathbb{E}(\mathbf{x} - \bar{x})(\mathbf{y} - \bar{y})^\top - r_{xy} R_y^{-1} \mathbb{E}(\mathbf{y} - \bar{y})(\mathbf{y} - \bar{y})^\top \\ &= r_{xy} - r_{xy} R_y^{-1} R_y \\ &= 0 \end{aligned} \quad (29.24)$$

But since $\mathbb{E} \tilde{\mathbf{x}} = 0$, we conclude that (29.23) holds. ■

It is because of the orthogonality property (29.23) between the estimation error

and the observation that equations (29.7) for w^o are referred to as the *normal equations*:

$$\boxed{R_y w^o = r_{yx}} \quad (\text{normal equations}) \quad (29.25)$$

It can be verified that these equations are always consistent, i.e., a solution w^o always exists independent of whether R_y is invertible or not — see Appendix 29.A. Moreover, the orthogonality condition (29.23) plays a critical role in characterizing linear estimators.

THEOREM 29.2. (Orthogonality principle) *An unbiased linear estimator is optimal in the least-mean-square-error sense if, and only if, its estimation error satisfies the orthogonality condition (29.23).*

Proof: One direction of the argument has already been proven prior to the statement, namely, the linear least-mean-square-error estimator, $\hat{\mathbf{x}}$, given by (29.8a), satisfies (29.23). With regards to the converse statement, assume now that we are given an unbiased linear estimator for \mathbf{x} of the form

$$\hat{\mathbf{x}}_u = \mathbf{y}^\top w_u - \theta_u \quad (29.26)$$

and that the corresponding estimation error satisfies the orthogonality condition (29.23), i.e., $\tilde{\mathbf{x}}_u \perp \mathbf{y}$. We verify that the parameters $\{w_u, \theta_u\}$ must coincide with the optimal parameters $\{w^o, \theta^o\}$ given by (29.7).

Indeed, the fact that $\hat{\mathbf{x}}_u$ is unbiased means that $\mathbb{E} \hat{\mathbf{x}}_u = \bar{x}$ and, hence, θ_u satisfies

$$\theta_u = -(\bar{x} - \bar{y}^\top w_u) \quad (29.27)$$

so that, by substituting into (29.26), we get that $\hat{\mathbf{x}}_u$ satisfies

$$\hat{\mathbf{x}}_u - \bar{x} = (\mathbf{y} - \bar{y})^\top w_u \quad (29.28)$$

Now using the assumed orthogonality condition $\tilde{\mathbf{x}}_u \perp \mathbf{y}$ we must have

$$\mathbb{E}(\mathbf{x} - \hat{\mathbf{x}}_u) \mathbf{y}^\top = 0 \quad (29.29)$$

which is equivalent to

$$\mathbb{E}(\mathbf{x} - \hat{\mathbf{x}}_u)(\mathbf{y} - \bar{y})^\top = 0 \quad (29.30)$$

since, by assumption, $\mathbb{E} \hat{\mathbf{x}}_u = \mathbb{E} \mathbf{x}$. Substituting expression (29.28) for $\hat{\mathbf{x}}_u$ into (29.30), we find that w_u must satisfy

$$\mathbb{E}(\mathbf{x} - \bar{x} - (\mathbf{y} - \bar{y})^\top w_u)(\mathbf{y} - \bar{y})^\top = 0 \quad (29.31)$$

which leads to

$$r_{yx} = R_y w_u \quad (29.32)$$

so that w^o and w_u satisfy the same normal equations. Substituting w_u by w^o into expression (29.27) we get that $\theta_u = \theta^o$. ■

Example 29.1 (Multiple noisy measurements of a binary signal) Consider a signal \mathbf{x} that assumes the values ± 1 with probability $1/2$ each. We collect N noisy measurements

$$\mathbf{y}_\ell = \mathbf{x} + \mathbf{v}_\ell, \quad \ell = 1, \dots, N \quad (29.33)$$

where \mathbf{v}_ℓ is zero-mean noise of unit-variance and independent of \mathbf{x} . We introduce the observation vector $\mathbf{y} = \text{col}\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$. Then, say, for $N = 5$, it is straightforward to find that

$$r_{xy} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}, \quad R_y = \begin{bmatrix} 2 & 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & 2 \end{bmatrix} \quad (29.34)$$

so that

$$\hat{\mathbf{x}} = r_{xy} R_y^{-1} \mathbf{y} = \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 & 2 \end{bmatrix}^{-1} \right) \mathbf{y} \quad (29.35)$$

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 . Observe that, for any N , the matrix R_y can be expressed as $R_y = I_N + \mathbf{1}\mathbf{1}^\top$, where I_N is the $N \times N$ identity matrix and $\mathbf{1}$ is the $N \times 1$ column vector with unit entries, $\mathbf{1} = \text{col}\{1, 1, 1, \dots, 1\}$. 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. 1.10). Specifically, it can be verified that, for any column vector $a \in \mathbb{R}^N$,

$$\left(I_N + aa^\top \right)^{-1} = I_N - \frac{aa^\top}{1 + \|a\|^2} \quad (29.36)$$

where $\|a\|^2 = a^\top a$ denotes the squared Euclidean norm of a . Using this result with $a = \mathbf{1}$, we find that

$$r_{xy} R_y^{-1} = \mathbf{1}^\top \left(I_N - \frac{\mathbf{1}\mathbf{1}^\top}{N+1} \right) = \mathbf{1}^\top - \frac{N}{N+1} \mathbf{1}^\top = \frac{\mathbf{1}^\top}{N+1} \quad (29.37)$$

so that

$$\hat{\mathbf{x}} = \frac{1}{N+1} \mathbf{1}^\top \mathbf{y} = \frac{1}{N+1} \sum_{\ell=1}^N \mathbf{y}_\ell \quad (29.38)$$

Recall that in this problem the variable \mathbf{x} is discrete and assumes the values ± 1 . The estimator $\hat{\mathbf{x}}$ in (29.38) will generally assume real-values. If one wishes to use $\hat{\mathbf{x}}$ to decide whether $\mathbf{x} = +1$ or $\mathbf{x} = -1$, then one may consider examining the sign of $\hat{\mathbf{x}}$ and use the sub-optimal estimator:

$$\hat{\mathbf{x}}_{\text{sub}} = \text{sign} \left(\frac{1}{N+1} \sum_{n=1}^N \mathbf{y}_n \right) \quad (29.39)$$

where the sign function was defined earlier in (27.36).

Example 29.2 (Learning a regression model from data) Assume we can estimate the price of a house in some neighborhood \mathbb{A} , measured in units of $\times 1000$ USD, in some affine manner from the surface area s (measured in m^2) and the unit's age a (measured in years), say, as:

$$\hat{P} = \alpha s + \beta a - \theta \quad (29.40)$$

for some unknown scalar parameters (α, β, θ) . Here, the scalar θ denotes an *offset* parameter and the above relation represents the equation of a plane mapping values

(s, a) into an estimate for the house price, denoted by \hat{P} . We collect the attributes $\{s, a\}$ into an observation vector:

$$y = \begin{bmatrix} s \\ a \end{bmatrix} \quad (29.41)$$

and the unknown parameters $\{\alpha, \beta\}$ into a column vector:

$$w = \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \quad (29.42)$$

Then, the mapping from y to \hat{P} can be written more compactly as:

$$\hat{P} = y^\top w - \theta \quad (29.43)$$

The price P plays the role of the variable \mathbf{x} that we wish to estimate from observations of \mathbf{y} . If we happen to know the first and second-order moments of the price and observation variables, then we could estimate (w, θ) by using

$$w^o = R_y^{-1} r_{yP}, \quad \theta^o = \bar{y}^\top w^o - \bar{P} \quad (29.44)$$

where \bar{P} is the average price of houses in neighborhood \mathbb{A} . Often, in practice, these statistical moments are not known beforehand. They can, however, be estimated from measurements. Assume we have available a list of N houses from neighborhood \mathbb{A} with their prices, size, and age; obviously, the house whose price we are interested in estimating should not be part of this list. We denote the available information by $\{P_n, y_n\}$ where $n = 1, 2, \dots, N$. Then, we can estimate the first and second-order moments from this data by using the sample averages:

$$\hat{P} = \frac{1}{N} \sum_{n=1}^N P_n, \quad \hat{y} = \frac{1}{N} \sum_{n=1}^N y_n \quad (29.45a)$$

and

$$\hat{r}_{yP} = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y})(P_n - \hat{P}) \quad (29.46a)$$

$$\hat{R}_y = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y})(y_n - \hat{y})^\top \quad (29.46b)$$

The parameters (w^o, θ^o) would be approximated by

$$w^* = \hat{R}_y^{-1} \hat{r}_{yP} \quad (29.47a)$$

$$\theta^* = \hat{y}^\top w^* - \hat{P} \quad (29.47b)$$

where we are using the star notation to refer to parameters estimated directly from data measurements; this will be a standard convention in our treatment.

Figure 29.1 illustrates these results by means of a simulation. The figure shows the scatter diagram for $N = 500$ points (P_n, y_n) representing the triplet (price, area, age) for a collection of 500 houses. The price is measured in units of $\times 1000$ USD, the area in units of m^2 , and the age in units of years. The spheres represent the measured data. The flat plane represents the regression plane that results from the above calculations, namely,

$$\hat{P} = y^\top w^* - \theta^* \quad (29.48)$$

with

$$w^* = \begin{bmatrix} 3.02 \\ -2.04 \end{bmatrix}, \quad \theta^* = 1.26 \quad (29.49)$$

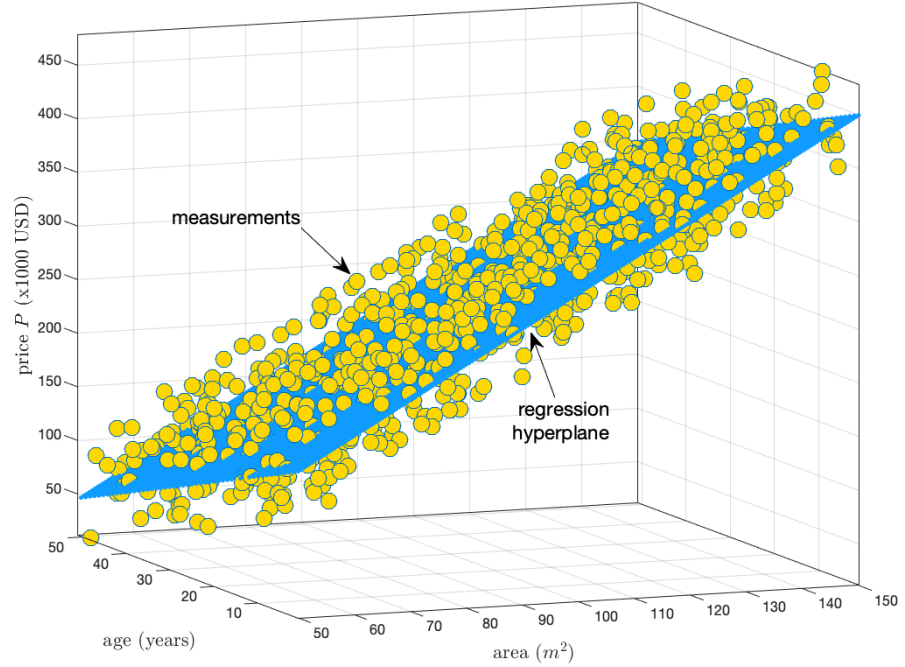


Figure 29.1 Scatter diagram of $N = 500$ points (P_n, y_n) representing the triplet (price, area, age) for a collection of 500 houses. The price is measured in units of $\times 1000$ USD, the area in units of m^2 , and the age in units of years. The spheres represent the measured data. The flat plane represents the fitted regression plane (29.48).

The values of the parameters (α, β, θ) are measured in units of $\times 1000$ USD. Now given a 17-year old house with area $102 m^2$, we can use the above parameter values to estimate/predict its price as follows:

$$\hat{P} = (3.02 \times 102) - (2.04 \times 17) - 1.26 = 272.1K \text{ USD} \quad (29.50)$$

29.2 CENTERING AND AUGMENTATION

We describe in this section two useful pre-processing steps that are common in inference and learning implementations in order to remove the need for the offset parameter, θ .

29.2.1 Centering

We start with centering. One useful fact to note is that the mean values $\{\bar{x}, \bar{y}\}$ appear subtracted from both sides of (29.8a). We say that the random variables

$\{\mathbf{x}, \mathbf{y}\}$ are being centered, by subtracting their means, and that the estimation problem actually amounts to estimating one centered variable from another centered variable, i.e., to estimating $\mathbf{x}_c = \mathbf{x} - \bar{x}$ from $\mathbf{y}_c = \mathbf{y} - \bar{y}$ in a linear manner. Indeed, note that the variables $\{\mathbf{x}, \mathbf{y}\}$ and $\{\mathbf{x}_c, \mathbf{y}_c\}$ have the same second-order moments since

$$R_{y_c} \triangleq \mathbb{E} \mathbf{y}_c \mathbf{y}_c^\top = \mathbb{E} (\mathbf{y} - \bar{y})(\mathbf{y} - \bar{y})^\top \triangleq R_y \quad (29.51a)$$

and

$$r_{x_c y_c} \triangleq \mathbb{E} \mathbf{x}_c (\mathbf{y}_c)^\top = \mathbb{E} (\mathbf{x} - \bar{x})(\mathbf{y} - \bar{y})^\top \triangleq r_{xy} \quad (29.51b)$$

so that estimating \mathbf{x}_c from \mathbf{y}_c leads to the same relation

$$\hat{\mathbf{x}}_c = r_{xy} R_y^{-1} \mathbf{y}_c \quad (29.52)$$

For this reason, and without loss in generality, it is customary in linear estimation problems to assume that the variables $\{\mathbf{x}, \mathbf{y}\}$ have zero means (or have already been centered) and to solve the estimation task under this condition.

The fact that centering arises in the context of linear least-mean-square-error estimation can also be seen directly from (29.5) if we were to impose the requirement that the estimator should be unbiased. In that case, we would conclude from (29.5) that θ must satisfy

$$\bar{x} = \bar{y}^\top w - \theta \implies \theta = \bar{y}^\top w - \bar{x} \quad (29.53)$$

Substituting this condition for θ into (29.5) we find that, in effect, the estimator we are seeking should be of the form

$$\hat{\mathbf{x}} - \bar{x} = (\mathbf{y} - \bar{y})^\top w \quad (29.54)$$

with centered variables appearing on both sides of the expression.

29.2.2 Augmentation

There is a second equivalent construction to centering, and which will be used extensively in later chapters. We refer to the affine model (29.5) and introduce the extended vectors:

$$\mathbf{y}' \triangleq \begin{bmatrix} 1 \\ \mathbf{y} \end{bmatrix}, \quad w' \triangleq \begin{bmatrix} -\theta \\ w \end{bmatrix}, \quad (M+1) \times 1 \quad (29.55)$$

where we have added the scalars 1 and $-\theta$ as leading entries on top of \mathbf{y} and w , respectively. Then, the affine estimator (29.5) can be rewritten in the *linear* (rather than affine) form:

$$\hat{\mathbf{x}} = (\mathbf{y}')^\top w' \quad (29.56)$$

and the design problem becomes one of determining a parameter vector $(w')^o$ that minimizes $\mathbb{E} \tilde{\mathbf{x}}^2$. We already know from the statement of Theorem 29.1 applied to this extended problem that $(w')^o$ should satisfy the linear equations:

$$R_{y'}(w')^o = r_{y'x} \quad (29.57)$$

where the first and second-order moments are computed as follows:

$$\bar{\mathbf{y}}' \triangleq \mathbb{E} \mathbf{y}' = \begin{bmatrix} 1 \\ \bar{y} \end{bmatrix} \quad (29.58a)$$

while

$$\begin{aligned} r_{y'x} &\triangleq \mathbb{E} (\mathbf{y}' - \bar{\mathbf{y}}')(\mathbf{x} - \bar{x}) \\ &= \mathbb{E} \begin{bmatrix} 0 \\ (\mathbf{y} - \bar{y}) \end{bmatrix} (\mathbf{x} - \bar{x}) \\ &= \begin{bmatrix} 0 \\ r_{yx} \end{bmatrix} \end{aligned} \quad (29.58b)$$

and

$$\begin{aligned} R_{y'} &\triangleq \mathbb{E} \mathbf{y}'(\mathbf{y}')^\top - \bar{\mathbf{y}}'(\bar{\mathbf{y}}')^\top \\ &= \mathbb{E} \begin{bmatrix} 1 & \mathbf{y}^\top \\ \mathbf{y} & \mathbf{y}\mathbf{y}^\top \end{bmatrix} - \begin{bmatrix} 1 & \bar{y}^\top \\ \bar{y} & \bar{y}\bar{y}^\top \end{bmatrix} \\ &= \begin{bmatrix} 0 & 0 \\ 0 & R_y \end{bmatrix} \end{aligned} \quad (29.58c)$$

If we denote the individual entries of $(w')^o$ by

$$(w')^o \triangleq \begin{bmatrix} -\theta^o \\ w^o \end{bmatrix} \quad (29.59)$$

then it follows from the expressions for $R_{y'}$ and $r_{y'x}$ and from (29.57) that the w^o component satisfies $R_y w^o = r_{yx}$. This result agrees with (29.7). Again, if we impose the unbiasedness condition that $\mathbb{E} \hat{\mathbf{x}} = \bar{x}$ at the optimal solution $(w')^o$, then we conclude from (29.56) that θ^o must satisfy

$$\begin{aligned} \bar{x} &= (\bar{\mathbf{y}}')^\top (w')^o \\ &= \begin{bmatrix} 1 & \bar{y}^\top \end{bmatrix} \begin{bmatrix} -\theta^o \\ w^o \end{bmatrix} \\ &= -\theta^o + \bar{y}^\top w^o \end{aligned} \quad (29.60)$$

from which we conclude that $\theta^o = \bar{y}^\top w^o - \bar{x}$, which again agrees with (29.7). For this reason, it is customary to assume that the variables have been extended to (\mathbf{y}', w') as in (29.55), and to seek a linear (as opposed to affine) model as in (29.56). It is a matter of convenience whether we assume that the variables (\mathbf{x}, \mathbf{y}) are centered and replaced by $(\mathbf{x}_c, \mathbf{y}_c)$ or that (\mathbf{y}, w) are extended and replaced by (\mathbf{y}', w') . The net effect in both cases is that we can assume that we are dealing with an *offset-free* problem that estimates \mathbf{x}_c from \mathbf{y}_c or \mathbf{x} from \mathbf{y}' in a linear (rather than affine) manner.

29.3 VECTOR ESTIMATION

The results from Sec. 29.1 can be easily extended to the case of estimating a *vector* (as opposed to a scalar) variable \mathbf{x} from multiple measurements. Thus, consider a column vector \mathbf{x} with entries

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_M \end{bmatrix}, \quad \bar{\mathbf{x}} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_M \end{bmatrix}, \quad (M \times 1) \quad (29.61)$$

We wish to estimate \mathbf{x} from multiple observations, say,

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}, \quad \bar{\mathbf{y}} = \begin{bmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \vdots \\ \bar{y}_N \end{bmatrix}, \quad (N \times 1) \quad (29.62)$$

The solution to this problem can be deduced from the results of the previous sections. First, according to (29.8a), we express the estimate for an arbitrary m -th entry of \mathbf{x} from all observations as follows:

$$(\hat{\mathbf{x}}_m - \bar{x}_m) = \mathbf{w}_m^T (\mathbf{y} - \bar{\mathbf{y}}) \quad (29.63)$$

for some column vector $\mathbf{w}_m \in \mathbb{R}^N$, and where the variables have been centered around their respective means. The vectors $\{\mathbf{w}_m\}$ are then determined by minimizing the aggregate mean-square-error:

$$\{\mathbf{w}_m^o\} \triangleq \underset{\{\mathbf{w}_m \in \mathbb{R}^N\}_{m=1}^M}{\operatorname{argmin}} \left\{ \sum_{m=1}^M \mathbb{E} (\mathbf{x}_m - \hat{\mathbf{x}}_m)^2 \right\} \quad (29.64)$$

29.3.1 Error Covariance Matrix

We can rewrite the cost in an equivalent form. Assume we collect the estimators into a column vector,

$$\hat{\mathbf{x}} = \begin{bmatrix} \hat{\mathbf{x}}_1 \\ \hat{\mathbf{x}}_2 \\ \vdots \\ \hat{\mathbf{x}}_M \end{bmatrix} \quad (29.65)$$

and introduce the error vector:

$$\tilde{\mathbf{x}} \triangleq \mathbf{x} - \hat{\mathbf{x}} \quad (29.66)$$

Then, the mean-squared Euclidean norm of $\tilde{\mathbf{x}}$ agrees with the cost appearing in (29.64), i.e.,

$$\mathbb{E} \|\tilde{\mathbf{x}}\|^2 = \sum_{m=1}^M \mathbb{E} (\mathbf{x}_m - \hat{\mathbf{x}}_m)^2 \quad (29.67)$$

If we further let $R_{\tilde{\mathbf{x}}}$ denote the covariance matrix of the error vector, namely,

$$R_{\tilde{\mathbf{x}}} \triangleq \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \quad (29.68)$$

then we also have that

$$\mathbb{E} \|\tilde{\mathbf{x}}\|^2 = \text{Tr}(R_{\tilde{\mathbf{x}}}) \quad (29.69)$$

so that, in effect, problem (29.64) is seeking the weight vectors $\{w_m^o\}$ that minimize the trace of the error covariance matrix:

$$\{w_m^o\} \triangleq \underset{\{w_m \in \mathbb{R}^N\}_{m=1}^M}{\text{argmin}} \left\{ \text{Tr}(R_{\tilde{\mathbf{x}}}) \right\} \quad (29.70)$$

We refer to the trace of $R_{\tilde{\mathbf{x}}}$, which appears in the above expression, as the mean-square-error in the vector case. We also refer to the matrix $R_{\tilde{\mathbf{x}}}$ as the mean-square-error matrix. In this way, for vector estimation problems, when we write m.s.e. we may be referring either to the scalar quantity $\text{Tr}(R_{\tilde{\mathbf{x}}})$ or to the matrix quantity $R_{\tilde{\mathbf{x}}}$ depending on the context. It is common to use the matrix representation for the m.s.e. in the vector case.

29.3.2 Normal Equations

Continuing with (29.64), we observe that the cost consists of the sum of M non-negative separable terms, with each term depending on the respective w_m . Therefore, we can determine the optimal coefficients $\{w_m, m = 1, \dots, M\}$ by minimizing each term separately:

$$w_m^o = \underset{w_m \in \mathbb{R}^N}{\text{argmin}} \mathbb{E} (\mathbf{x}_m - \hat{\mathbf{x}}_m)^2 \quad (29.71)$$

This is the same problem we studied before: estimating a scalar variable \mathbf{x}_m from multiple observations $\{\mathbf{y}_\ell, \ell = 1, \dots, N\}$. We already know that the solution is given by the normal equations:

$$R_y w_m^o = r_{y\mathbf{x}_m}, \quad m = 1, 2, \dots, M \quad (29.72)$$

where $r_{\mathbf{x}_m \mathbf{y}}$ is the cross-covariance vector of \mathbf{x}_m with \mathbf{y} :

$$r_{\mathbf{x}_m \mathbf{y}} \triangleq \mathbb{E} (\mathbf{x}_m - \bar{\mathbf{x}}_m)(\mathbf{y} - \bar{\mathbf{y}})^\top \quad (29.73)$$

Moreover, the resulting estimation error satisfies the orthogonality condition:

$$\mathbb{E} \tilde{\mathbf{x}}_m \mathbf{y}^\top = 0, \quad m = 1, \dots, M \quad (29.74)$$

where $\tilde{\mathbf{x}}_m = \mathbf{x}_m - \hat{\mathbf{x}}_m$. Note that the $\{r_{x_my}\}$ are the rows of the $M \times N$ cross-covariance matrix R_{xy} between the vectors \mathbf{x} and \mathbf{y} :

$$R_{xy} = \mathbb{E}(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{y} - \bar{\mathbf{y}})^\top = \begin{bmatrix} r_{x_1y} \\ r_{x_2y} \\ \vdots \\ r_{x_My} \end{bmatrix}, \quad (M \times N) \quad (29.75)$$

Therefore, by collecting the w_m^o from (29.72) as columns into a matrix W^o , for all $k = 1, 2, \dots, M$:

$$W^o = \begin{bmatrix} w_1^o & w_2^o & \dots & w_M^o \end{bmatrix}, \quad (N \times M) \quad (29.76)$$

and by noting that $R_{yx} = R_{xy}^\top$, we find that W^o satisfies the normal equations

$$\boxed{R_y W^o = R_{yx}} \quad (\text{normal equations}) \quad (29.77)$$

It follows that the optimal estimator is given by

$$\hat{\mathbf{x}} - \bar{\mathbf{x}} = (W^o)^\top (\mathbf{y} - \bar{\mathbf{y}}) \quad (29.78)$$

or, equivalently,

$$\boxed{\hat{\mathbf{x}} - \bar{\mathbf{x}} = R_{xy} R_y^{-1} (\mathbf{y} - \bar{\mathbf{y}})} \quad (29.79)$$

In view of (29.74), the resulting estimation error vector satisfies the orthogonality condition:

$$\boxed{\mathbb{E} \tilde{\mathbf{x}} \mathbf{y}^\top = 0} \quad (\text{orthogonality principle}) \quad (29.80)$$

where $\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$. The corresponding minimum mean-square error matrix is given by

$$\begin{aligned} \text{m.m.s.e.} &= \mathbb{E}(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^\top \\ &= \mathbb{E}(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \hat{\mathbf{x}})^\top \\ &= \mathbb{E}(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^\top \\ &= \mathbb{E}((\mathbf{x} - \bar{\mathbf{x}}) - (\hat{\mathbf{x}} - \bar{\mathbf{x}}))(\mathbf{x} - \bar{\mathbf{x}})^\top \\ &= \mathbb{E}(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^\top - \mathbb{E}(\hat{\mathbf{x}} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^\top \end{aligned} \quad (29.81)$$

and we conclude from (29.79) that

$$\boxed{\text{m.m.s.e.} = R_x - R_{xy} R_y^{-1} R_{yx}} \quad (29.82)$$

in terms of the covariance and cross-covariance matrices of \mathbf{x} and \mathbf{y} .

29.4 LINEAR MODELS

We apply the mean-square-error estimation theory of the previous sections to the important case of linear models, which arises often in applications. Thus, assume that zero-mean random vectors $\{\mathbf{x}, \mathbf{y}\}$ are related via a linear model of the form:

$$\mathbf{y} = H\mathbf{x} + \mathbf{v} \quad (29.83)$$

for some $N \times M$ known matrix H . We continue to assume that \mathbf{y} is $N \times 1$ and \mathbf{x} is $M \times 1$ so that we are estimating a vector variable from another vector variable. We explained earlier that the zero-mean assumption is not restrictive since the random variables \mathbf{x} and \mathbf{y} can be assumed to have been centered. In the above model, the variable \mathbf{v} denotes a zero-mean random additive noise vector with known covariance matrix, $R_v = \mathbb{E}\mathbf{v}\mathbf{v}^\top$. The covariance matrix of \mathbf{x} is also assumed to be known, say, $R_x = \mathbb{E}\mathbf{x}\mathbf{x}^\top$. Both $\{\mathbf{x}, \mathbf{v}\}$ are uncorrelated with each other, i.e., $\mathbb{E}\mathbf{x}\mathbf{v}^\top = 0$, and we further assume that

$$R_x > 0, \quad R_v > 0 \quad (29.84)$$

Two equivalent representations

According to (29.79), when $R_y > 0$, the linear least-mean-square-error estimator of \mathbf{x} given \mathbf{y} is

$$\hat{\mathbf{x}} = R_{xy}R_y^{-1}\mathbf{y} \quad (29.85)$$

Because of (29.83), the covariances $\{R_{xy}, R_y\}$ can be determined in terms of the linear model parameters $\{H, R_x, R_v\}$. Indeed, the uncorrelatedness of $\{\mathbf{x}, \mathbf{v}\}$ gives

$$R_y = \mathbb{E}\mathbf{y}\mathbf{y}^\top = \mathbb{E}(H\mathbf{x} + \mathbf{v})(H\mathbf{x} + \mathbf{v})^\top = HR_xH^\top + R_v \quad (29.86)$$

$$R_{xy} = \mathbb{E}\mathbf{x}\mathbf{y}^\top = \mathbb{E}\mathbf{x}(H\mathbf{x} + \mathbf{v})^\top = R_xH^\top \quad (29.87)$$

Moreover, since $R_v > 0$ we immediately get $R_y > 0$. Expression (29.85) for $\hat{\mathbf{x}}$ then becomes

$$\hat{\mathbf{x}} = R_xH^\top (R_v + HR_xH^\top)^{-1}\mathbf{y} \quad (29.88)$$

This expression can be rewritten in an equivalent form by using the *matrix inversion formula* (1.81). The result states that for arbitrary matrices $\{A, B, C, D\}$ of compatible dimensions, if A and C are invertible, then

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \quad (29.89)$$

Applying this identity to the matrix $(R_v + HR_xH^\top)^{-1}$ in (29.88), with the identifications

$$A \leftarrow R_v, \quad B \leftarrow H, \quad C \leftarrow R_x, \quad D \leftarrow H^\top \quad (29.90)$$

we obtain

$$\begin{aligned}
 \hat{\mathbf{x}} &= R_x H^\top \left\{ R_v^{-1} - R_v^{-1} H (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \right\} \mathbf{y} \\
 &= \left\{ R_x (R_x^{-1} + H^\top R_v^{-1} H) - R_x H^\top R_v^{-1} H \right\} (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \mathbf{y} \\
 &= (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \mathbf{y}
 \end{aligned} \tag{29.91}$$

where in the second equality we factored out $(R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \mathbf{y}$ from the right. Hence,

$$\hat{\mathbf{x}} = (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \mathbf{y} \tag{29.92}$$

This alternative form is useful in several contexts. Observe, for example, that when H happens to be a column vector (i.e., when \mathbf{x} is a scalar), the quantity $(R_v + H R_x H^\top)$ that appears in (29.88) is a matrix, while the quantity $(R_x^{-1} + H^\top R_v^{-1} H)$ that appears in (29.92) is a scalar. In this case, the representation (29.92) leads to a simpler expression for $\hat{\mathbf{x}}$. In general, the convenience of using (29.88) or (29.92) depends on the situation at hand.

It further follows that the m.m.s.e. matrix is given by

$$\begin{aligned}
 \text{m.m.s.e.} &= \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top = \mathbb{E} \mathbf{x} \mathbf{x}^\top - \mathbb{E} \hat{\mathbf{x}} \hat{\mathbf{x}}^\top \\
 &= \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}}) (\mathbf{x} - \hat{\mathbf{x}})^\top, \quad \text{since } \tilde{\mathbf{x}} \perp \hat{\mathbf{x}} \\
 &= R_x - (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} H R_x \\
 &= (R_x^{-1} + H^\top R_v^{-1} H)^{-1}
 \end{aligned} \tag{29.93}$$

where in the last equality we used the matrix inversion lemma again. That is,

$$\text{m.m.s.e.} = (R_x^{-1} + H^\top R_v^{-1} H)^{-1} \tag{29.94}$$

LEMMA 29.1. (Equivalent linear estimators) Consider two zero-mean random vectors $\{\mathbf{x}, \mathbf{y}\}$ that are related via a linear model of the form $\mathbf{y} = H\mathbf{x} + \mathbf{v}$. The variables $\{\mathbf{x}, \mathbf{v}\}$ have zero mean, positive-definite covariance matrices R_x and R_v , respectively, and are uncorrelated with each other. The linear least-mean-square-error estimator of \mathbf{x} given \mathbf{y} can be computed by either of the following equivalent expressions:

$$\hat{\mathbf{x}} = R_x H^\top (R_v + H R_x H^\top)^{-1} \mathbf{y} \tag{29.95a}$$

$$= (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \mathbf{y} \tag{29.95b}$$

with the resulting m.m.s.e. value given by

$$\text{m.m.s.e.} = (R_x^{-1} + H^\top R_v^{-1} H)^{-1} \tag{29.96}$$

29.5 DATA FUSION

We illustrate one application of the theory for linear models to the important problem of fusing information from several sources in order to enhance the accuracy of the estimation process. Thus, assume that we have a collection of N sensors that are distributed over some region in space. All sensors are interested in estimating the same zero-mean vector \mathbf{x} with covariance matrix $R_x = \mathbb{E} \mathbf{x} \mathbf{x}^\top$. For example, the sensors could be tracking a moving object and their objective is to estimate the speed and direction of motion of the target.

Assume each sensor k collects a measurement vector \mathbf{y}_k that is related to \mathbf{x} via a linear model, say,

$$\mathbf{y}_k = H_k \mathbf{x} + \mathbf{v}_k \quad (29.97)$$

where H_k is the model matrix that maps \mathbf{x} to \mathbf{y}_k at sensor k and \mathbf{v}_k is zero-mean measurement noise with covariance matrix $R_k = \mathbb{E} \mathbf{v}_k \mathbf{v}_k^\top$. In general, the quantity \mathbf{y}_k is a vector, say, of size $L_k \times 1$. If the vector \mathbf{x} has size $M \times 1$, then H_k is $L_k \times M$. We assume $L_k \geq M$ so that each sensor k has at least as many measurements as the size of the unknown vector \mathbf{x} .

29.5.1 Fusing Raw Data

In the first fusion method, each sensor k transmits its measurement vector \mathbf{y}_k and its model parameters $\{H_k, R_k, R_x\}$ to a remote fusion center. The latter collects all measurements from across the nodes, $\{\mathbf{y}_k, k = 1, 2, \dots, N\}$, and all model parameters $\{(H_k, R_k), k = 1, 2, \dots, N\}$. The collected data satisfies the model:

$$\underbrace{\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} H_1 \\ H_2 \\ \vdots \\ H_N \end{bmatrix}}_H \mathbf{x} + \underbrace{\begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_N \end{bmatrix}}_{\mathbf{v}} \quad (29.98)$$

which is a linear model of the same form as (29.83). We assume that the noises $\{\mathbf{v}_k\}$ across all nodes are uncorrelated with each other so that the covariance matrix of the aggregate noise vector, \mathbf{v} , is block diagonal:

$$R_v = \text{blkdiag}\{R_1, R_2, \dots, R_N\} \quad (29.99)$$

Now, using (29.92), the fusion center can determine the estimator of \mathbf{x} that is based on all measurements $\{\mathbf{y}_k\}$ as follows:

$$\hat{\mathbf{x}} = (R_x^{-1} + H^\top R_v^{-1} H)^{-1} H^\top R_v^{-1} \mathbf{y} \quad (29.100)$$

The resulting m.m.s.e. would be

$$P = (R_x^{-1} + H^\top R_v^{-1} H)^{-1} \quad (29.101)$$

where we are denoting the m.m.s.e. by the letter P ; it is a matrix of size $M \times M$. We note that we are deliberately using form (29.92) for the estimator of \mathbf{x} because, as we are going to see shortly, this form will allow us to derive a second more efficient fusion method.

The solution method (29.100) requires all sensors to transmit $\{\mathbf{y}_k, H_k, R_k\}$ to the fusion center; this amounts to a total of

$$L_k + L_k M + \frac{1}{2} L_k^2 \quad \text{entries to be transmitted per sensor} \quad (29.102)$$

where the last term ($L_k^2/2$) arises from the transmission of (half of) the entries of the $L_k \times L_k$ symmetric matrix R_k .

29.5.2 Fusing Processed Data

In the fusion method (29.100), the fusion center fuses the *raw* data $\{\mathbf{y}_k, H_k, R_k\}$ that are collected at the sensors. A more efficient data fusion method is possible and leads to a reduction in the amount of communication resources that are necessary between the sensors and the fusion center. This alternative method is based on the sensors performing some local processing first and then sharing the results of the processing step with the fusion center. The two fusion modes are illustrated in Fig. 29.2.

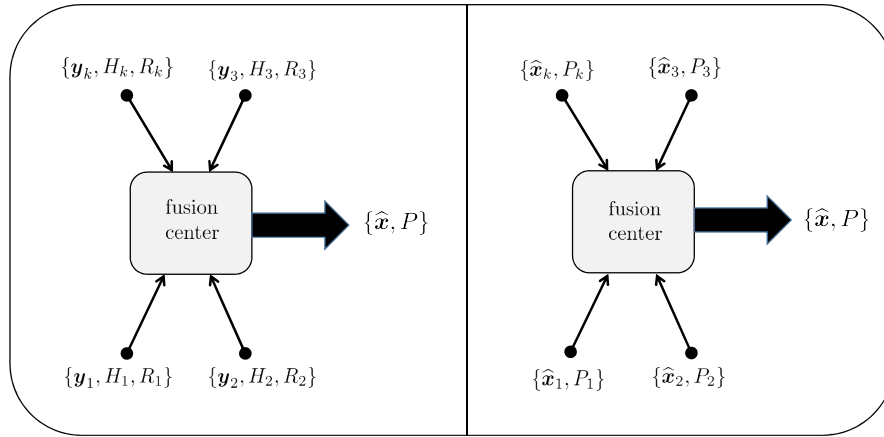


Figure 29.2 Two modes for data fusion: (a) on the left, the sensors share their raw data $\{\mathbf{y}_k, H_k, R_k\}$ with the fusion center, (b) while on the right, the sensors share their local estimation results $\{\hat{\mathbf{x}}_k, P_k\}$ with the fusion center.

Specifically, assume that each node estimates \mathbf{x} using its own data \mathbf{y}_k . We denote the resulting estimator by $\hat{\mathbf{x}}_k$ and it is given by:

$$\hat{\mathbf{x}}_k = (R_x^{-1} + H_k^T R_k^{-1} H_k)^{-1} H_k^T R_k^{-1} \mathbf{y}_k \quad (29.103)$$

The corresponding m.m.s.e. is

$$P_k = (R_x^{-1} + H_k^\top R_k^{-1} H_k)^{-1} \quad (29.104)$$

We assume now that the nodes share the processed data $\{\hat{\mathbf{x}}_k, P_k\}$ with the fusion center rather than the raw data $\{\mathbf{y}_k, H_k, R_k\}$. It turns out that the desired global quantities $\{\hat{\mathbf{x}}, P\}$ can be recovered from these processed data.

To begin with, observe that we can rework expression (29.101) for the global m.m.s.e. as follows:

$$\begin{aligned} P^{-1} &= R_x^{-1} + H^\top R_v^{-1} H \\ &= R_x^{-1} + \sum_{k=1}^N H_k^\top R_k^{-1} H_k \\ &= \sum_{k=1}^N (R_x^{-1} + H_k^\top R_k^{-1} H_k) - (N-1)R_x^{-1} \\ &= \sum_{k=1}^N P_k^{-1} - (N-1)R_x^{-1} \end{aligned} \quad (29.105)$$

This expression allows the fusion center to determine P^{-1} directly from knowledge of the quantities $\{P_k^{-1}, R_x^{-1}\}$. Note further that the global expression (29.100) can be rewritten as

$$P^{-1}\hat{\mathbf{x}} = H^\top R_v^{-1} \mathbf{y} \quad (29.106)$$

which, using H and R_v from (29.98)–(29.99), leads to

$$P^{-1}\hat{\mathbf{x}} = \sum_{k=1}^N H_k^\top R_k^{-1} \mathbf{y}_k = \sum_{k=1}^N P_k^{-1} \hat{\mathbf{x}}_k \quad (29.107)$$

Therefore, we arrive at the following conclusion to fuse the data from multiple sensors.

LEMMA 29.2. (Data fusion) *Consider a collection of N linear measurements of the form $\mathbf{y}_k = H_k \mathbf{x} + \mathbf{v}_k$, where \mathbf{x} and \mathbf{v}_k have zero mean, positive-definite covariance matrices R_x and R_k , respectively, and are uncorrelated with each other. Let $\hat{\mathbf{x}}$ denote the linear least-mean-square-error estimator of \mathbf{x} given the N observations $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$ with error covariance matrix $P = \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top$. Let $\hat{\mathbf{x}}_k$ denote the linear least-mean-square-error estimator of the same \mathbf{x} given only \mathbf{y}_k with error covariance matrix $P_k = \mathbb{E} \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top$. It holds that*

$$P^{-1} = (P_1^{-1} + P_2^{-1} + \dots + P_N^{-1}) - (N-1)R_x^{-1} \quad (29.108a)$$

$$P^{-1}\hat{\mathbf{x}} = P_1^{-1}\hat{\mathbf{x}}_1 + P_2^{-1}\hat{\mathbf{x}}_2 + \dots + P_N^{-1}\hat{\mathbf{x}}_N \quad (29.108b)$$

Observe from (29.108b) that the individual estimators are scaled by the inverses of their m.m.s.e matrices so that more accurate estimators are given more weight.

In this method, the sensors need to send to the fusion center the processed information $\{\hat{\mathbf{x}}_k, P_k\}$. This amounts to a total of

$$M + M^2/2 \text{ entries to be transmitted per node} \quad (29.109)$$

which is smaller than (29.102) given that $L_k \geq M$.

29.6 MINIMUM-VARIANCE UNBIASED ESTIMATION

In the previous section we examined the linear model (29.83) where the unknown, \mathbf{x} , is modeled as a *random* variable with covariance matrix, R_x . We encountered one instance of this model earlier in Example 29.1, which dealt with the problem of estimating a zero-mean scalar random variable, \mathbf{x} , from a collection of noisy measurements, $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$. The model relating the variables in that example is a special case of (29.83) since it amounts to

$$\underbrace{\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}}_{\triangleq \mathbf{y}} = \underbrace{\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}}_{\triangleq H} \mathbf{x} + \underbrace{\begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_N \end{bmatrix}}_{\triangleq \mathbf{v}} \quad (29.110)$$

where

$$H = \text{col}\{1, 1, \dots, 1\} \quad (29.111a)$$

$$R_x = \sigma_x^2 \quad (29.111b)$$

$$R_v = \sigma_v^2 I_N \quad (29.111c)$$

Note that, for generality, we are using generic values for σ_x^2 and σ_v^2 rather than set them equal to one, as was the case in Example 29.1. We can recover the solution (29.38) by appealing to expression (29.92), which gives

$$\hat{\mathbf{x}} = \left(\frac{1}{\sigma_x^2} + \frac{N}{\sigma_v^2} \right)^{-1} \frac{1}{\sigma_v^2} \sum_{\ell=1}^N \mathbf{y}_\ell = \frac{1}{N + \frac{1}{\text{SNR}}} \sum_{\ell=1}^N \mathbf{y}_\ell \quad (29.112)$$

in terms of the signal-to-noise ratio defined as

$$\text{SNR} \triangleq \sigma_x^2 / \sigma_v^2 \quad (29.113)$$

In (29.110), the variable \mathbf{x} is assumed to have been selected at random and, subsequently, N noisy measurements of this same value are collected. The observations are used to estimate \mathbf{x} according to (29.112). Observe that the solution does *not* correspond to computing the sample mean of the observations. Expression (29.112) would reduce to the sample mean only when $\text{SNR} \rightarrow \infty$.

We now consider an alternative formulation for estimating the unknown by modeling it as a *deterministic* unknown constant, say, x , rather than a *random*

quantity, \mathbf{x} . For this purpose, we replace the earlier linear model (29.83) by one of the form

$$\mathbf{y} = Hx + \mathbf{v} \quad (29.114)$$

where, compared with (29.83), we are replacing the boldface letter \mathbf{x} by the normal letter x (remember that we reserve the boldface notation to random variables). The observation vector \mathbf{y} in (29.114) continues to be random since the disturbance \mathbf{v} is random. Any estimator for x that is based on \mathbf{y} will also be a random variable itself. Given model (29.114), we now study the problem of designing an optimal linear estimator for x of the form

$$\hat{x} = W^\top \mathbf{y} \quad (29.115)$$

for some coefficient matrix $W^\top \in \mathbb{R}^{M \times N}$ to be determined. It will turn out that, for such problems, W is found by solving a *constrained* least-mean-square-error estimation problem, as opposed to the unconstrained estimation problem (29.70).

29.6.1 Problem Formulation

Thus, consider a zero-mean random noise variable \mathbf{v} with a positive-definite covariance matrix $R_v = \mathbb{E} \mathbf{v} \mathbf{v}^\top > 0$, and let \mathbf{y} be a noisy measurement of Hx according to model (29.114) where x is the unknown constant vector that we wish to estimate. The dimensions of the data matrix H are denoted by $N \times M$ and it is assumed that $N \geq M$ and that H has full rank:

$$\text{rank}(H) = M, \quad N \geq M \quad (29.116)$$

That is, H is a tall matrix so that the number of available measurements is at least as many as the number of unknown entries in x . The full rank condition on H guarantees that the matrix product $H^\top R_v^{-1} H$ is positive-definite — recall result (1.59). The inverse of this matrix product will appear in the expression for the estimator.

We are interested in determining a linear estimator for x of the form $\hat{x} = W^\top \mathbf{y}$. The choice of W should satisfy two conditions:

(a) (Unbiasedness). That is, we must guarantee $\mathbb{E} \hat{x} = x$, which is the same as $W^\top \mathbb{E} \mathbf{y} = x$. But from (29.114) we have $\mathbb{E} \mathbf{y} = Hx$ so that W must satisfy $W^\top Hx = x$, no matter what the value of x is. This condition means that W should satisfy

$$W^\top H = I_M \quad (29.117)$$

(b) (Optimality). The choice of W should minimize the trace of the covariance matrix of the estimation error, $\tilde{x} = x - \hat{x}$. Using the condition $W^\top H = I_M$, we find that

$$\hat{x} = W^\top \mathbf{y} = W^\top (Hx + \mathbf{v}) = W^\top Hx + W^\top \mathbf{v} = x + W^\top \mathbf{v} \quad (29.118)$$

so that $\tilde{\mathbf{x}} = -W\mathbf{v}$. This means that the error covariance matrix, as a function of W , is given by

$$\mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T = W^T R_v W \quad (29.119)$$

Combining (29.117) and (29.119), we conclude that the desired W can be found by solving the following constrained optimization problem:

$$W^o \triangleq \underset{W}{\operatorname{argmin}} \left\{ \operatorname{Tr}(W^T R_v W) \right\}, \text{ subject to } W^T H = I_M, \quad R_v > 0 \quad (29.120)$$

The estimator $\hat{\mathbf{x}} = (W^o)^T \mathbf{y}$ that results from the solution of (29.120) is known as the *minimum-variance-unbiased estimator*, or MVUE for short. It is also sometimes called the *best linear unbiased estimator* or BLUE.

Example 29.3 (Guessing the solution) Let us first try to guess the form of the solution to the constrained problem (29.120) by appealing to the solution of the linear least-mean-square-error estimation problem (29.83). In that formulation, the unknown, \mathbf{x} , is modeled as a random variable with covariance matrix R_x . From expression (29.92), we know that the linear estimator is given by

$$\hat{\mathbf{x}} = (R_x^{-1} + H^T R_v^{-1} H)^{-1} H^T R_v^{-1} \mathbf{y} \quad (29.121)$$

Now assume that the covariance matrix of \mathbf{x} has the particular form $R_x = \sigma_x^2 I$, with a sufficiently large positive scalar σ_x^2 (i.e., $\sigma_x^2 \rightarrow \infty$). That is, assume that the variance of each of the entries of \mathbf{x} is “infinitely” large. In this way, the variable \mathbf{x} can be “regarded” as playing the role of some unknown constant vector, x . Then, the above expression for $\hat{\mathbf{x}}$ reduces to

$$\hat{\mathbf{x}} = (H^T R_v^{-1} H)^{-1} H^T R_v^{-1} \mathbf{y} \quad (29.122)$$

This conclusion suggests that the choice

$$(W^o)^T = (H^T R_v^{-1} H)^{-1} H^T R_v^{-1} \quad (29.123)$$

should solve the problem of estimating the unknown constant vector x for model (29.114). We establish this result more formally next.

29.6.2 Gauss-Markov Theorem

Result (29.123) is a manifestation of the Gauss-Markov Theorem.

THEOREM 29.3. (Gauss-Markov theorem) Consider a linear model of the form $\mathbf{y} = H\mathbf{x} + \mathbf{v}$, where x is an unknown constant, \mathbf{v} has zero-mean and covariance matrix $R_v > 0$, and H is a tall full-rank matrix (with as least as many rows as columns). The minimum-variance unbiased estimator for x , the one that solves (29.120), is given by

$$\hat{\mathbf{x}}_{\text{MVUE}} = (H^T R_v^{-1} H)^{-1} H^T R_v^{-1} \mathbf{y} \quad (29.124)$$

Equivalently, the optimal W in (29.120) is

$$(W^o)^T = (H^T R_v^{-1} H)^{-1} H^T R_v^{-1} \quad (29.125)$$

Moreover, the resulting minimum mean-square error is

$$R_{\hat{\mathbf{x}}} \triangleq \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T = (H^T R_v^{-1} H)^{-1} \quad (29.126)$$

Proof: Let $\mathcal{J}(W) = W^T R_v W$ denote the cost function that appears in (29.120). Some straightforward algebra shows that $\mathcal{J}(W)$ can be expressed as

$$\mathcal{J}(W) = (W - W^o)^T R_v (W - W^o) + (W^o)^T R_v W^o \quad (29.127)$$

This is because, using $W^T H = I$,

$$\begin{aligned} W^T R_v W^o &= W^T R_v \left(R_v^{-1} H (H^T R_v^{-1} H)^{-1} \right) \\ &= W^T H (H^T R_v^{-1} H)^{-1} = (H^T R_v^{-1} H)^{-1} \end{aligned} \quad (29.128)$$

Likewise, $(W^o)^T R_v W^o = (H^T R_v^{-1} H)^{-1}$. Relation (29.127) expresses the cost $\mathcal{J}(K)$ as the sum of two nonnegative-definite terms: one is independent of W and is equal to $(W^o)^T R_v W^o$, while the other is dependent on W . It is then clear, since $R_v > 0$, that the trace of the cost is minimized by choosing $W = W^o$. Note further that the matrix W^o satisfies the constraint $(W^o)^T H = I_M$. ■

Example 29.4 (Sample mean estimator) Let us reconsider problem (29.110) where \mathbf{x} is now modeled as an unknown constant, i.e., we now write

$$\underbrace{\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}}_{\triangleq \mathbf{y}} = \underbrace{\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}}_{\triangleq H} x + \underbrace{\begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_N \end{bmatrix}}_{\triangleq \mathbf{v}} \quad (29.129)$$

where the boldface \mathbf{x} is replaced by x . In this case, the value of x can be regarded as the mean value for each measurement, \mathbf{y}_ℓ . Using expression (29.124) we find that the minimum-variance unbiased estimator for x is given by (compare with (29.112)):

$$\hat{\mathbf{x}}_{\text{MVUE}} = \frac{1}{N} \sum_{\ell=1}^N \mathbf{y}_\ell \quad (29.130)$$

which is simply the sample mean estimator that we are familiar with from introductory courses on statistics.

29.7 COMMENTARIES AND DISCUSSION

Linear estimation. In this chapter we covered the basics of linear mean-square error regression analysis and highlighted only concepts that are relevant to the subject matter of the book, motivated by the presentations in Sayed (2003,2008) and Kailath, Sayed, and Hassibi (2000). The pioneering work in this domain was done independently by the Russian mathematician **Andrey Kolmogorov (1903–1987)** in the works by Kolmogorov (1939,1941a,b) and the American mathematician **Norbert Wiener (1894–1964)** in the work by Wiener (1949); the latter reference was originally published in 1942 as a classified report during World War II. Kolmogorov was motivated by the work of Wold (1938) on stationary processes and solved a linear prediction problem for discrete-time stationary random processes. Wiener, on the other hand, solved a continuous-time prediction problem under causality constraints by means of an elegant technique now known as the Wiener-Hopf technique introduced in Wiener and Hopf (1931). Readers interested in more details about Wiener's contribution, and linear estimation theory in general, may consult the textbook by Kailath, Sayed, and Hassibi (2000).

Unbiased estimators. In the least-mean-square-error estimation problems studied here, the estimators were required to be unbiased. Sometimes, unbiasedness can be a hurdle to minimizing the mean-square error. This is because there are estimators that are biased but that can achieve smaller error variances than unbiased estimators — see, e.g., Rao (1973), Cox and Hinkley (1974), and Kendall and Stuart (1976–1979).

Two interesting examples to this effect are the following given in Sayed (2008): the first example is from Kay (1993, pp. 310–311) while the second example is from Rao (1973). In the derivation leading to (29.130) we studied the problem of estimating the mean value, x , of N measurements $\{\mathbf{y}_\ell\}$. The minimum-variance unbiased estimator for x was seen to be given by the sample mean estimator:

$$\hat{\mathbf{x}}_{\text{MVUE}} = \frac{1}{N} \sum_{\ell=1}^N \mathbf{y}_\ell \quad (29.131)$$

The value of x was not restricted in any way; it was only assumed to be an unknown constant and that it could assume any value in the interval $(-\infty, \infty)$. But what if we know beforehand that x is limited to some interval, say $[-\alpha, \alpha]$ for some finite $\alpha > 0$? One way to incorporate this piece of information into the design of an estimator for x is to consider the following alternative construction:

$$\tilde{\mathbf{x}} = \begin{cases} -\alpha, & \text{if } \hat{\mathbf{x}}_{\text{MVUE}} < -\alpha \\ \hat{\mathbf{x}}_{\text{MVUE}}, & \text{if } -\alpha \leq \hat{\mathbf{x}}_{\text{MVUE}} \leq \alpha \\ \alpha, & \text{if } \hat{\mathbf{x}}_{\text{MVUE}} > \alpha \end{cases} \quad (29.132)$$

in terms of a realization for $\hat{\mathbf{x}}_{\text{MVUE}}$. In this way, $\tilde{\mathbf{x}}$ will always assume values within $[-\alpha, \alpha]$. A calculation in Kay (1993) shows that although the above (truncated mean) estimator $\tilde{\mathbf{x}}$ is biased, it nevertheless satisfies $\mathbb{E}(x - \tilde{\mathbf{x}})^2 < \mathbb{E}(x - \hat{\mathbf{x}}_{\text{MVUE}})^2$ — see Prob. 29.5. In other words, the truncated mean estimator results in a smaller mean-square error.

A second classical example from the realm of statistics is the variance estimator. In this case, the parameter to be estimated is the variance of a random variable \mathbf{y} given access to several observations of it, say $\{\mathbf{y}_\ell\}$. Let σ_y^2 denote the variance of \mathbf{y} . Two well-known estimators for σ_y^2 are

$$\hat{\sigma}_y^2 = \frac{1}{N-1} \sum_{\ell=1}^N (\mathbf{y}_\ell - \bar{\mathbf{y}})^2 \quad \text{and} \quad \widetilde{\sigma}_y^2 = \frac{1}{N+1} \sum_{\ell=1}^N (\mathbf{y}_\ell - \bar{\mathbf{y}})^2 \quad (29.133)$$

where $\bar{y} = \frac{1}{N} \sum_{\ell=1}^N \mathbf{y}_\ell$ is the sample mean. The first one is unbiased while the second one is biased. However, it is shown in Rao (1973) that

$$\mathbb{E} \left(\sigma_y^2 - \widetilde{\sigma}_y^2 \right)^2 < \mathbb{E} \left(\sigma_y^2 - \widehat{\sigma}_y^2 \right)^2 \quad (29.134)$$

We therefore see that biased estimators can result in smaller mean-square errors. However, unbiasedness is often a desirable property in practice since it guarantees that, on average, the estimator agrees with the unknown quantity that we seek to estimate.

Gauss-Markov theorem. Theorem 29.3 characterizes unbiased linear estimators of smallest error variance (or covariance), also known as BLUE estimators. Given an unknown $x \in \mathbb{R}^M$ and a random observation $\mathbf{y} \in \mathbb{R}^N$, the theorem was obtained by solving the constrained optimization problem (29.120), namely,

$$\min_W \mathbb{E} \|\tilde{\mathbf{x}}\|^2, \text{ subject to } \hat{\mathbf{x}} = W^\top \mathbf{y}, \mathbf{y} = Hx + \mathbf{v}, W^\top H = I_M, \mathbb{E} \mathbf{v} \mathbf{v}^\top = R_v \quad (29.135)$$

The significance of the solution (29.124) is perhaps best understood if we recast it in the context of least-squares problems. We will study such problems in greater detail in Chapter 50; see also Sayed (2008). For now, let us assume that $R_v = \sigma_v^2 I_N$, so that the noise components are uncorrelated with each other and have equal variances. Then, expression (29.124) for the estimate of x reduces to

$$\hat{x} = (H^\top H)^{-1} H^\top y \quad (29.136)$$

We are going to see later that this expression can be interpreted as the solution to the following least-squares problem. Given a noisy *deterministic* observation vector y satisfying

$$y = Hx + v \quad (29.137)$$

the unknown x can be estimated by solving

$$\hat{x} = \operatorname{argmin}_{x \in \mathbb{R}^M} \|y - Hx\|^2 \quad (29.138)$$

in terms of the squared Euclidean norm of the difference $y - Hx$. Indeed, if we expand the squared error we get

$$J(x) \triangleq \|y - Hx\|^2 = \|y\|^2 - 2y^\top Hx + x^\top H^\top Hx \quad (29.139)$$

Setting the gradient vector relative to x to zero at \hat{x} gives

$$\nabla_x J(x) \Big|_{x=\hat{x}} = 2\hat{x}^\top H^\top H - 2y^\top H = 0 \implies \hat{x} = (H^\top H)^{-1} H^\top y \quad (29.140)$$

which is the same expression we had in (29.136). We therefore find that the standard least-squares estimate for x coincides with the BLUE estimate. More generally, consider a *weighted* least-squares problem of the form

$$\hat{x} = \operatorname{argmin}_{x \in \mathbb{R}^M} (y - Hx)^\top R (y - Hx) \quad (29.141)$$

where $R > 0$ denotes some weighting matrix. Differentiating again and solving for \hat{x} gives

$$\hat{x} = (H^\top R H)^{-1} H^\top R y \quad (29.142)$$

Comparing with (29.124) we find that the weighted least-squares solution would agree with the BLUE estimate if we select the weighting matrix as $R = R_v^{-1}$ (i.e., as the inverse of the noise covariance matrix). Therefore, the Gauss-Markov theorem is essentially stating that the least-squares solution leads to the best linear unbiased estimator (BLUE) if the weighting matrix is matched with R_v^{-1} .

The original version of the Gauss–Markov theorem with $R_v = \sigma_v^2 I_N$ is due to the German mathematician **Carl Friedrich Gauss (1777–1855)** and the Russian mathematician **Andrey Markov (1856–1922)**, who published versions of the result in Gauss (1821) and Markov (1912). The extension to the case of arbitrary covariance matrices R_v was given by Aitken (1935) — see also the overview by Plackett (1949,1950). The discussion in Sec. 29.6 on the Gauss Markov theorem, and the leading Secs. 29.4 and 29.5 on linear models and data fusion are adapted from the discussion in Kailath, Sayed, and Hassibi (2000).

PROBLEMS¹

29.1 Show that the linear least-mean-square-error estimator defined by (29.70) also minimizes the determinant of the error covariance matrix, $\det(R_{\tilde{x}})$.

29.2 Show that the linear least-mean-square-error estimator defined by (29.70) also minimizes the $\mathbb{E} \tilde{\mathbf{x}}^T W \tilde{\mathbf{x}}$ for any $W \geq 0$.

29.3 All variables are zero mean. Show that for any three random variables $\{\mathbf{x}, \mathbf{y}, \mathbf{z}\}$ it holds that

$$\hat{\mathbf{x}}_{y,z} = \hat{\mathbf{x}}_y + (\hat{\mathbf{x}}_y)_{\tilde{z}_y}$$

where

$$\left\{ \begin{array}{ll} \hat{\mathbf{x}}_{y,z} &= \text{linear least-mean-squares estimator (l.l.m.s.e) of } \mathbf{x} \text{ given } \{\mathbf{z}, \mathbf{y}\}. \\ \hat{\mathbf{x}}_y &= \text{l.l.m.s.e of } \mathbf{x} \text{ given } \mathbf{y}. \\ \hat{\mathbf{z}}_y &= \text{l.l.m.s.e of } \mathbf{z} \text{ given } \mathbf{y}. \\ \tilde{\mathbf{x}}_y &= \mathbf{x} - \hat{\mathbf{x}}_y \\ \tilde{\mathbf{z}}_y &= \mathbf{z} - \hat{\mathbf{z}}_y \\ (\tilde{\mathbf{x}}_y)_{\tilde{z}_y} &= \text{l.l.m.s.e of } \tilde{\mathbf{x}}_y \text{ given } \tilde{\mathbf{z}}_y. \end{array} \right.$$

What is the geometric interpretation of this result?

29.4 Verify that the mean-square-error values that correspond to the estimators $\{\hat{\mathbf{x}}^\circ, \hat{\mathbf{x}}^\bullet\}$ defined by (29.149a)–(29.149b) coincide.

29.5 Refer to the truncated mean estimator (29.132). Show that it results in a smaller mean-square error, namely, $\mathbb{E}(x - \tilde{\mathbf{x}})^2 < \mathbb{E}(x - \hat{\mathbf{x}}_{\text{MVUE}})^2$.

29.6 Let $\{\mathbf{x}, \mathbf{y}\}$ denote two zero-mean random variables with positive-definite covariance matrices $\{R_x, R_y\}$. Let $\hat{\mathbf{x}}$ denote the linear least-mean-square-error estimator of \mathbf{x} given \mathbf{y} . Likewise, let $\hat{\mathbf{y}}$ denote the linear least-mean-square-error estimator of \mathbf{y} given \mathbf{x} . Introduce the estimation errors $\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$ and $\tilde{\mathbf{y}} = \mathbf{y} - \hat{\mathbf{y}}$, and denote their covariance matrices by $R_{\tilde{x}}$ and $R_{\tilde{y}}$, respectively.

(a) Show that $R_x R_{\tilde{x}}^{-1} \hat{\mathbf{x}} = R_{xy} R_{\tilde{y}}^{-1} \mathbf{y}$.

(b) Assume $\{\mathbf{y}, \mathbf{x}\}$ are related via a linear model of the form $\mathbf{y} = H\mathbf{x} + \mathbf{v}$, where H is a matrix of appropriate dimensions while \mathbf{v} has zero-mean with covariance matrix R_v and is uncorrelated with \mathbf{x} . Verify that the identity of part (a) reduces to $R_{\tilde{x}}^{-1} \hat{\mathbf{x}} = H^T R_v^{-1} \mathbf{y}$.

29.7 Let \mathbf{x} be a zero-mean random variable with an $M \times M$ positive-definite covariance matrix R_x . Let $\hat{\mathbf{x}}_1$ denote the linear least-mean-square-error estimator of \mathbf{x} given a zero-mean observation \mathbf{y}_1 . Likewise, let $\hat{\mathbf{x}}_2$ denote the linear least-mean-square-error estimator of the same variable \mathbf{x} given a second zero-mean observation \mathbf{y}_2 . That is, we have two separate estimators for \mathbf{x} from two separate sources. Let P_1 and P_2 denote the corresponding error covariance matrices: $P_1 = \mathbb{E} \tilde{\mathbf{x}}_1 \tilde{\mathbf{x}}_1^T$ and $P_2 = \mathbb{E} \tilde{\mathbf{x}}_2 \tilde{\mathbf{x}}_2^T$ where

¹ Some problems in this section are adapted from exercises in Sayed (2003,2008) and Kailath, Sayed, and Hassibi (2000).

$\tilde{\mathbf{x}}_j = \mathbf{x} - \hat{\mathbf{x}}_j$, Assume $P_1 > 0$ and $P_2 > 0$ and that the cross-covariance matrix

$$\mathbb{E} \begin{bmatrix} \mathbf{x} \\ \mathbf{y}_1 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y}_2 \end{bmatrix}^\top$$

has rank M .

- (a) Show that the linear least-mean-square-error estimator of \mathbf{x} given *both* $\{\mathbf{y}_1, \mathbf{y}_2\}$, denoted by $\hat{\mathbf{x}}$, satisfies $P^{-1}\hat{\mathbf{x}} = P_1^{-1}\hat{\mathbf{x}}_1 + P_2^{-1}\hat{\mathbf{x}}_2$, where P denotes the resulting error covariance matrix and is given by $P^{-1} = P_1^{-1} + P_2^{-1} - R_x^{-1}$.
- (b) Assume $\{\mathbf{y}_1, \mathbf{x}\}$ and $\{\mathbf{y}_2, \mathbf{x}\}$ are related via linear models of the form $\mathbf{y}_1 = H_1\mathbf{x} + \mathbf{v}_1$ and $\mathbf{y}_2 = H_2\mathbf{x} + \mathbf{v}_2$, where $\{\mathbf{v}_1, \mathbf{v}_2\}$ have zero means with covariance matrices $\{R_{v_1}, R_{v_2}\}$ and are uncorrelated with each other and with \mathbf{x} . Verify that this situation satisfies the required rank-deficiency condition and conclude that the estimator of \mathbf{x} given $\{\mathbf{y}_1, \mathbf{y}_2\}$ is given by the expression in part (a).

29.8 Let $\mathbf{y}_1 = H_1\mathbf{x} + \mathbf{v}_1$ and $\mathbf{y}_2 = H_2\mathbf{x} + \mathbf{v}_2$ denote two linear observation models with the same unknown random vector \mathbf{x} . All random variables have zero-mean. The covariance and cross-covariance matrices of $\{\mathbf{x}, \mathbf{v}_1, \mathbf{v}_2\}$ are denoted by

$$\mathbb{E} \begin{bmatrix} \mathbf{x} \\ \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}^\top = \begin{bmatrix} R_x & 0 & 0 \\ 0 & R_1 & C \\ 0 & C^\top & R_2 \end{bmatrix}$$

In particular, observe that we are assuming the noises to be correlated with $C = \mathbb{E}\mathbf{v}_1\mathbf{v}_2^\top$. All covariance matrices are assumed to be invertible whenever necessary.

- (a) Show how you would replace the observation vectors $\{\mathbf{y}_1, \mathbf{y}_2\}$ by two other observation vectors $\{\mathbf{z}_1, \mathbf{z}_2\}$ of similar dimensions such that they satisfy linear models of the form

$$\mathbf{z}_1 = G_1\mathbf{x} + \mathbf{w}_1, \quad \mathbf{z}_2 = G_2\mathbf{x} + \mathbf{w}_2$$

for some matrices G_1 and G_2 to be specified, and where the noises $\{\mathbf{w}_1, \mathbf{w}_2\}$ are now uncorrelated. What are the covariance matrices of \mathbf{w}_1 and \mathbf{w}_2 in terms of R_1 and R_2 ?

- (b) Let $\hat{\mathbf{x}}_1$ be the linear least-mean-square-error estimator (l.l.m.s.e.) of \mathbf{x} given \mathbf{z}_1 with error covariance matrix P_1 . Similarly, let $\hat{\mathbf{x}}_2$ be the l.l.m.s.e. of \mathbf{x} given \mathbf{y}_2 with error covariance matrix P_2 . Let further $\hat{\mathbf{x}}$ denote the l.l.m.s.e. of \mathbf{x} given $\{\mathbf{y}_1, \mathbf{y}_2\}$ with error covariance matrix P . Determine expressions for $\hat{\mathbf{x}}$ and P in terms of $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, P_1, P_2, C, R_x, R_1, R_2\}$.

29.9 Let $\mathbf{y} = H\mathbf{x} + \mathbf{v}$. All random variables have zero-mean. The covariance and cross-covariance matrices of $\{\mathbf{x}, \mathbf{v}\}$ are denoted by

$$\mathbb{E} \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix}^\top = \begin{bmatrix} R_x & C \\ C^\top & R_v \end{bmatrix}$$

with positive-definite R_x and R_v .

- (a) What is the l.l.m.s.e. of \mathbf{x} given \mathbf{y} ? What is the corresponding m.m.s.e.?
- (b) A new scalar observation, α , is added to \mathbf{y} and a new row vector is added to H such that

$$\begin{bmatrix} \mathbf{y} \\ \alpha \end{bmatrix} = \begin{bmatrix} H \\ h^\top \end{bmatrix} \mathbf{x} + \begin{bmatrix} \mathbf{v} \\ n \end{bmatrix}$$

where h^\top is a row vector, n is uncorrelated with all other variables and has variance σ^2 . Let $\hat{\mathbf{x}}_{\text{new}}$ denote the new estimator for \mathbf{x} given $\{\mathbf{y}, \alpha\}$. Relate $\hat{\mathbf{x}}_{\text{new}}$ to $\hat{\mathbf{x}}$ from part (a). Relate also their m.m.s.e. values.

29.10 All variables are zero-mean. Let

$$\begin{bmatrix} \mathbf{y}_a \\ \mathbf{y} \\ \mathbf{y}_b \end{bmatrix} = \begin{bmatrix} H_a \\ H \\ H_b \end{bmatrix} \mathbf{x} + \begin{bmatrix} \mathbf{v}_a \\ \mathbf{v} \\ \mathbf{v}_b \end{bmatrix}$$

where $\{\mathbf{v}_a, \mathbf{v}, \mathbf{v}_b\}$, are uncorrelated with \mathbf{x} and have zero mean and covariance matrices:

$$\mathbb{E} \begin{bmatrix} \mathbf{v}_a \\ \mathbf{v} \\ \mathbf{v}_b \end{bmatrix} \begin{bmatrix} \mathbf{v}_a \\ \mathbf{v} \\ \mathbf{v}_b \end{bmatrix}^\top = \begin{bmatrix} R_a & S_a & 0 \\ S_a^\top & R & S_b \\ 0 & S_b^\top & R_b \end{bmatrix}$$

Let $\hat{\mathbf{x}}_{\mathbf{y}_a, \mathbf{y}}$ denote the linear estimator of \mathbf{x} given $\{\mathbf{y}_a, \mathbf{y}\}$. Let $\hat{\mathbf{x}}_{\mathbf{y}_b, \mathbf{y}}$ denote the linear estimator of \mathbf{x} given $\{\mathbf{y}_b, \mathbf{y}\}$. Can you relate these estimators, and their minimum-mean-square-error, to each other?

29.11 Let $\mathbf{y} = \mathbf{x} + \mathbf{v}$, where \mathbf{x} and \mathbf{v} are independent zero-mean Gaussian random variables with variances σ_x^2 and σ_v^2 , respectively. Show that the linear least-mean-square-error estimator of \mathbf{x}^2 using $\{\mathbf{y}, \mathbf{y}^2\}$ is

$$\widehat{\mathbf{x}^2} = \sigma_x^2 + \frac{\sigma_x^4}{\sigma_x^4 + 2\sigma_x^2\sigma_v^2 + \sigma_v^4} (\mathbf{y}^2 - \sigma_x^2 - \sigma_v^2)$$

29.12 A random variable \mathbf{z} is defined as follows

$$\mathbf{z} = \begin{cases} -\mathbf{x}, & \text{with probability } p \\ H\mathbf{x} + \mathbf{v}, & \text{with probability } 1 - p \end{cases}$$

where \mathbf{x} and \mathbf{v} are zero-mean uncorrelated random vectors. Assume we know the linear least-mean-square-error estimator of \mathbf{x} given \mathbf{y} , namely, $\hat{\mathbf{x}}_{|\mathbf{y}}$, where \mathbf{y} is a zero-mean random variable that is also uncorrelated with \mathbf{v} .

- Find an expression for $\hat{\mathbf{z}}_{|\mathbf{y}}$ in terms of $\hat{\mathbf{x}}_{|\mathbf{y}}$.
- Find an expression for the linear least-mean-square-error estimator $\hat{\mathbf{x}}_{|\mathbf{z}}$ and the corresponding m.m.s.e.

29.13 Consider the distributed network with m nodes, shown in Fig. 29.13. Each node k observes a zero-mean measurement \mathbf{y}_k that is related to an unknown zero-mean variable \mathbf{x} via a linear model of the form $\mathbf{y}_k = H_k \mathbf{x} + \mathbf{v}_k$, where the data matrix H_k is known, and the noise \mathbf{v}_k is zero mean and uncorrelated with \mathbf{x} . The noises across all nodes are uncorrelated with each other. Let $\{R_x, R_k\}$ denote the positive-definite covariance matrices of $\{\mathbf{x}, \mathbf{v}_k\}$, respectively. Introduce the following notation:

- At each node k , the notation $\hat{\mathbf{x}}_k$ denotes the linear least-mean-squares estimator of \mathbf{x} that is based on the observation \mathbf{y}_k . Likewise, P_k denotes the resulting error covariance matrix, $P_k = \mathbb{E} \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top$.
- At each node k , the notation $\hat{\mathbf{x}}_{1:k}$ denotes the linear least-mean-squares estimator of \mathbf{x} that is based on the observations $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k\}$, i.e., on the observations collected at nodes 1 through k . Likewise, $P_{1:k}$ denotes the resulting error covariance matrix, $P_{1:k} = \mathbb{E} \tilde{\mathbf{x}}_{1:k} \tilde{\mathbf{x}}_{1:k}^\top$.

The network functions as follows. Node 1 uses \mathbf{y}_1 to estimate \mathbf{x} . The resulting estimator, $\hat{\mathbf{x}}_1$, and the corresponding error covariance matrix, $P_1 = \mathbb{E} \tilde{\mathbf{x}}_1 \tilde{\mathbf{x}}_1^\top$, are transmitted to node 2. Node 2 in turn uses its measurement \mathbf{y}_2 and the data $\{\hat{\mathbf{x}}_1, P_1\}$ received from node 1 to compute the estimator of \mathbf{x} that is based on both observations $\{\mathbf{y}_1, \mathbf{y}_2\}$. Note that node 2 *does not* have access to \mathbf{y}_1 but only to \mathbf{y}_2 and the information received from node 1. The estimator computed by node 2, $\hat{\mathbf{x}}_{1:2}$, and the corresponding error covariance matrix, $P_{1:2}$, are then transmitted to node 3. Node 3 evaluates $\{\hat{\mathbf{x}}_{1:3}, P_{1:3}\}$ using $\{\mathbf{y}_3, \hat{\mathbf{x}}_{1:2}, P_{1:2}\}$ and transmits $\{\hat{\mathbf{x}}_{1:3}, P_{1:3}\}$ to node 4 and so forth.

- Find an expression for $\hat{\mathbf{x}}_{1:m}$ in terms of $\hat{\mathbf{x}}_{1:m-1}$ and $\hat{\mathbf{x}}_m$.
- Find an expression for $P_{1:m}^{-1}$ in terms of $\{P_{1:m-1}^{-1}, P_m^{-1}, R_x^{-1}\}$.
- Find a recursion relating $P_{1:m}$ to $P_{1:m-1}$.
- Show that $P_{1:m}$ is a non-increasing sequence as a function of m .

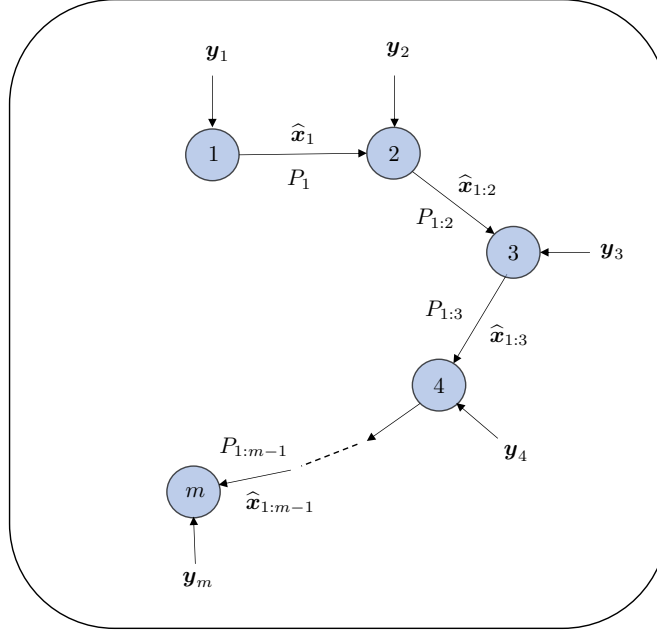


Figure 29.3 A distributed network with m nodes for Prob. 29.13.

- (e) Assume $H_k = H$ for all k and $R_{v_k} = R_v > 0$. Assume further that H is tall and has full column rank. Find $\lim_{m \rightarrow \infty} P_{1:m}$.

29.14 Consider two sensors labeled $k = 1, 2$, and assume each sensor has an unbiased estimator, $\{\mathbf{w}_k, k = 1, 2\}$ for some $M \times 1$ column vector \mathbf{w}^o . Let $\{P_k, k = 1, 2\}$ denote the error covariance matrix, $P_k = \mathbb{E}(\mathbf{w}^o - \mathbf{w}_k)(\mathbf{w}^o - \mathbf{w}_k)^\top$. Assume the errors of the two estimators are uncorrelated, i.e., $\mathbb{E}(\mathbf{w}^o - \mathbf{w}_1)(\mathbf{w}^o - \mathbf{w}_2)^\top = 0$. Consider a new aggregate estimator of the form $\hat{\mathbf{w}} = \alpha \mathbf{w}_1 + (1 - \alpha) \mathbf{w}_2$.

- (a) If α is nonnegative, determine the optimal scalar α that minimizes the mean-square-error, i.e., $\min_{\alpha \geq 0} \mathbb{E} \|\mathbf{w}^o - \hat{\mathbf{w}}\|^2$.
 (b) Repeat part (a) when α is not restricted to being nonnegative. When would a negative α be advantageous?
 (c) Let $P = \mathbb{E}(\mathbf{w}^o - \hat{\mathbf{w}})(\mathbf{w}^o - \hat{\mathbf{w}})^\top$. How does P compare to P_1 and P_2 in both cases (a) and (b)?
 (d) Now assume the errors of the two estimators are *correlated* instead, i.e., $\mathbb{E}(\mathbf{w}^o - \mathbf{w}_1)(\mathbf{w}^o - \mathbf{w}_2)^\top = C$, for some matrix C . Repeat parts (a)–(c).

29.15 Consider N sensors labeled $k = 1, 2, \dots, N$. Each node has an unbiased estimate of some unknown column vector $\mathbf{w}^o \in \mathbb{R}^M$. We denote the individual estimator at node k by \mathbf{w}_k . We also denote the error covariance matrix of \mathbf{w}_k by P_k and the cross-covariance matrix of \mathbf{w}_k and \mathbf{w}_ℓ by $P_{k\ell}$. A sensor S wishes to combine the estimators $\{\mathbf{w}_k, k = 1, \dots, N\}$ through $\hat{\mathbf{w}}_S = \sum_{k=1}^N a_k \mathbf{w}_k$ in order to optimize the cost function:

$$\min_{\{a_k\}} \mathbb{E} \left\| \mathbf{w}^o - \sum_{k=1}^N a_k \mathbf{w}_k \right\|^2$$

where the $\{a_k\}$ are real-valued scalars.

- (a) Find a condition on the coefficients $\{a_k\}$ to ensure that the resulting \hat{w}_S is an unbiased estimator for w^o .
- (b) Under condition (a), find the optimal coefficients $\{a_k\}$. Your solution should not depend on w^o .
- (c) Assume the reliability of each estimator w_k is measured by the scalar $\sigma_k^2 = \text{Tr}(P_k)$. The smaller the σ_k^2 is, the more reliable the estimator will be. What is the relation between the optimal coefficients $\{a_k\}$ and the reliability factors $\{\sigma_k^2\}$?
- (d) Evaluate the reliability of the estimator \hat{w}_S .
- (e) Motivate and derive a stochastic gradient algorithm for updating the coefficients $\{a_k\}$ in part (b).
- (f) How is the estimator of part (b) different from the unbiased linear least-mean-squares estimator of w^o based on the $\{w_k\}$? Find the latter estimator.
- (g) Find the minimum mean-square-error (m.m.s.e.) of the estimators in parts (b) and (f) for the case where $P_{k\ell} = 0$ when $\ell \neq k$. Specialize your result to the case $P_k = P$ for all k and compare the resulting mean-square-errors.

29.16 Consider a collection of N independent and identically-distributed random variables, $\{\mathbf{y}(n), n = 0, 1, \dots, N-1\}$. Each $\mathbf{y}(n)$ has a Gaussian distribution with zero mean and variance σ_y^2 . We want to use the observations $\{\mathbf{y}(n)\}$ to estimate the variance σ_y^2 in the following manner:

$$\hat{\sigma}_y^2 = \alpha \sum_{n=0}^{N-1} \mathbf{y}^2(n)$$

for some scalar parameter α to be determined.

- (a) What is the mean of the estimator $\hat{\sigma}_y^2$ in terms of α and σ_y^2 ?
- (b) Evaluate the mean-square-error below in terms of α and σ_y^2 :

$$\text{m.s.e.} = \mathbb{E}(\hat{\sigma}_y^2 - \sigma_y^2)^2$$

- (c) Determine the optimal scalar α that minimizes the m.s.e.. Is the corresponding estimator biased or unbiased?
- (d) For what value of α would the estimator be unbiased? What is the m.s.e. of this estimator and how does it compare to the m.s.e. of the estimator from part (c)?

29.17 Consider noisy observations $\mathbf{y}(n) = \mathbf{x} + \mathbf{v}(n)$, where \mathbf{x} and $\mathbf{v}(n)$ are independent random variables, $\mathbf{v}(n)$ is a white random process with zero mean and distributed as follows:

$\mathbf{v}(n)$ is Gaussian with variance σ_v^2 with probability q
 $\mathbf{v}(n)$ is uniformly distributed over $[-a, a]$ with probability $1 - q$

Moreover, \mathbf{x} assumes the values ± 1 with equal probability. The value of \mathbf{x} is the same for all measurements $\{\mathbf{y}(n)\}$. All variables are real-valued.

- (a) Find an expression for the linear least-mean-square-error estimator (l.l.m.s.e.) of \mathbf{x} given the collection of N observations $\{\mathbf{y}(0), \mathbf{y}(1), \dots, \mathbf{y}(N-1)\}$.
- (b) Find the l.l.m.s.e. of \mathbf{x} given the observations $\{\mathbf{y}(0), \mathbf{y}(1), \dots, \mathbf{y}(N-1)\}$ and $\{\mathbf{y}^2(0), \mathbf{y}^2(1), \dots, \mathbf{y}^2(N-1)\}$. How does the answer compare to part (a)?

29.18 This problem deals with constrained mean-square-error estimation. Let \mathbf{d} denote a scalar zero-mean random variable with variance σ_d^2 , and let \mathbf{u} denote an $M \times 1$ zero-mean random vector with covariance matrix $R_u = \mathbb{E} \mathbf{u} \mathbf{u}^T > 0$. Consider the constrained optimization problem

$$\min_w \mathbb{E}(\mathbf{d} - \mathbf{u}^T w)^2, \quad \text{subject to } c^T w = \alpha$$

where c is a known $M \times 1$ vector and α is a known real scalar.

- (a) Let $z = w - R_u^{-1} r_{ud}$ and $r_{du} = \mathbb{E} \mathbf{d} \mathbf{u}^T$. Show that the above optimization problem is equivalent to the following:

$$\min_{z \in \mathbb{R}^M} \left\{ \sigma_d^2 - r_{du} R_u^{-1} r_{ud} + z^T R_u z \right\}, \quad \text{subject to } c^T z = \alpha - c^T R_u^{-1} r_{ud}$$

- (b) Show that the optimal solution, w^o , of the constrained optimization problem is given by

$$w^o = R_u^{-1} r_{ud} - \left(\frac{c^T R_u^{-1} r_{ud} - \alpha}{c^T R_u^{-1} c} \right) R_u^{-1} c$$

Verify that this solution satisfies the constraint $c^T w^o = \alpha$.

29.19 Let \mathbf{x} be a zero-mean random variable with an $M \times M$ positive-definite covariance matrix R_x . Let $\hat{\mathbf{x}}_{\mathbf{y}_1}$ denote the linear least-mean-square-error estimator of \mathbf{x} given a zero-mean observation \mathbf{y}_1 with covariance matrix R_{y_1} . Likewise, let $\hat{\mathbf{x}}_{\mathbf{y}_2}$ denote the linear least-mean-square-error estimator of \mathbf{x} given another zero-mean observation \mathbf{y}_2 with covariance matrix R_{y_2} . Let $R_{y_1, y_2} = \mathbb{E} \mathbf{y}_1 \mathbf{y}_2^T$. We want to determine another estimator for \mathbf{x} by combining $\hat{\mathbf{x}}_{\mathbf{y}_1}$ and $\hat{\mathbf{x}}_{\mathbf{y}_2}$ in a convex manner as follows:

$$\hat{\mathbf{x}} = \lambda \hat{\mathbf{x}}_{\mathbf{y}_1} + (1 - \lambda) \hat{\mathbf{x}}_{\mathbf{y}_2}$$

where λ is a real scalar lying inside the interval $0 \leq \lambda \leq 1$.

- (a) Determine the value of λ that results in an estimator $\hat{\mathbf{x}}$ with the smallest mean-square error.
 (b) If λ is allowed to be any arbitrary real scalar (not necessarily limited to the range $[0, 1]$), how much smaller can the mean-square-error be?

29.20 Let $\mathbf{y} = \mathbf{s} + \mathbf{v}$ be a vector of measurements, where \mathbf{v} is noise and \mathbf{s} is the desired signal. Both \mathbf{v} and \mathbf{s} are zero-mean uncorrelated random vectors with covariance matrices $\{R_v, R_s\}$, respectively. We wish to determine a unit-norm column vector, w , such that the signal-to-noise ratio in the output signal, $\mathbf{y}^T w$, is maximized.

- (a) Verify that the covariance matrices of the signal and noise components in $\mathbf{y}^T w$ are equal to $w^T R_s w$ and $w^T R_v w$, respectively.
 (b) Assume first that $R_v = \sigma_v^2 I$. Use the Rayleigh-Ritz characterization (1.16) to conclude that the solution of

$$\max_{\|w\|=1} \left(\frac{w^T R_s w}{\sigma_v^2 \|w\|^2} \right)$$

is given by the unit-norm eigenvector that corresponds to the maximum eigenvalue of R_s , written as $w^o = q_{\max}$, where $R_s q_{\max} = \lambda_{\max} q_{\max}$. Verify further that the resulting maximum SNR is equal to $\lambda_{\max} / \sigma_v^2$.

- (c) Assume now that \mathbf{v} is colored noise so that its covariance matrix is not necessarily diagonal. Introduce the eigen-decomposition $R_v = U \Lambda U^T$, where U is orthogonal and Λ is diagonal with positive entries. Let $L = U \Lambda^{1/2}$. Repeat the argument of part (b) to show that the solution of

$$\max_{\|w\|=1} \left(\frac{w^T R_s w}{w^T R_v w} \right)$$

is now related to the unit-norm eigenvector that corresponds to the maximum eigenvalue of $L^{-1} R_s (L^T)^{-1}$.

29.21 Consider the optimization problem:

$$\min_W W^T R_v W, \quad \text{subject to } W^T H = A, \quad R_v > 0$$

where W^T is $M \times N$, H is $N \times P$, A is $M \times P$, $P < N$, $M < N$, and H has full rank. In the text we assumed A is square and equal to the identity matrix (see (29.120)). Show that the optimal solution is given by

$$(W^o)^T = A(H^T R_v^{-1} H)^{-1} H^T R_v^{-1}$$

and that the resulting minimum cost is $A(H^T R_v^{-1} H)^{-1} A^T$.

29.22 Refer to (29.108a). Compare P to P_k , for each $k = 1, 2, \dots, N$. Specifically, verify that the difference $P_k - P$ is non-negative definite.

29.23 Refer to (29.108a)–(29.108b). Let $\{\hat{\mathbf{x}}, P\}$ be the estimator and the m.m.s.e. that result from estimating \mathbf{x} from data across all N sensors. Let $\{\hat{\mathbf{x}}', P'\}$ be the estimator and the m.m.s.e. that result from estimating \mathbf{x} from data across the first $N - 1$ sensors. Relate $\{\hat{\mathbf{x}}, P\}$ to $\{\hat{\mathbf{x}}', P'\}$.

29.24 All variables are zero-mean. Consider a complex-valued scalar random variable \mathbf{d} and a complex-valued $M \times 1$ regression vector \mathbf{u} . Let

$$\hat{\mathbf{d}} = (w^o)^\top \mathbf{u} = \mathbf{u}^\top w^o$$

denote the linear least-mean-square-error (l.l.m.s.e.) estimator of \mathbf{d} given \mathbf{u} for some $M \times 1$ vector w^o . Consider additionally the problem of estimating separately the real and imaginary parts of \mathbf{d} using knowledge of the real and imaginary parts of \mathbf{u} , also in the linear least-mean-square-error sense, namely,

$$\hat{\mathbf{d}}_{\text{real}} = (w_{\text{real}}^o)^\top \begin{bmatrix} \text{Re}(\mathbf{u}) \\ \text{Im}(\mathbf{u}) \end{bmatrix}, \quad \hat{\mathbf{d}}_{\text{imag}} = (w_{\text{imag}}^o)^\top \begin{bmatrix} \text{Re}(\mathbf{u}) \\ \text{Im}(\mathbf{u}) \end{bmatrix}$$

for some $2M \times 1$ vectors w_{real}^o and w_{imag}^o .

- Argue that estimating the real and imaginary parts of \mathbf{d} from the real and imaginary parts of \mathbf{u} is equivalent to estimating the real and imaginary parts of \mathbf{d} from $\{\mathbf{u}, \mathbf{u}^*\}$, where \mathbf{u}^* is the complex conjugate transpose of \mathbf{u} .
- What are the optimal choices for w^o , w_{real}^o and w_{imag}^o ?
- Let $\hat{\mathbf{d}}_2 = \hat{\mathbf{d}}_{\text{real}} + j\hat{\mathbf{d}}_{\text{imag}}$ denote the estimator that is obtained for \mathbf{d} from this second construction. What is the corresponding m.m.s.e.? How does it compare to the m.m.s.e. obtained for $\hat{\mathbf{d}} = (w^o)^\top \mathbf{u}$? Under what conditions will both constructions lead to the same m.m.s.e.?

29.A CONSISTENCY OF NORMAL EQUATIONS

In this appendix we verify that the normal equations (29.25) are always consistent, i.e., we establish that a solution w^o always exists. Moreover, the solution is either unique or there are infinitely many solutions. In the latter case, all solutions will differ by vectors in the nullspace of R_y and, moreover, all of them will lead to the *same* estimator for \mathbf{x} and to the *same* mean-square error. Only these possibilities can occur.

Proof (Consistency of normal equations). We verify that at least one solution w^o exists to the normal equations $R_y w^o = r_{yx}$. For this purpose, we need to verify that r_{yx} belongs to the range space of R_y , i.e.,

$$r_{yx} \in \mathcal{R}(R_y) \quad (29.143)$$

We show this property by contradiction. Assume that (29.143) does not hold. Under this assumption, there should exist some nonzero vector $p \in \mathcal{N}(R_y)$ that is *not* orthogonal to r_{yx} , namely,

$$\exists p \text{ such that } R_y p = 0, \quad r_{yx}^\top p \neq 0 \quad (29.144)$$

It follows from $R_y p = 0$ that $p^\top R_y p = 0$ so that

$$\begin{aligned} p^\top R_y p = 0 &\iff p^\top \left(\mathbb{E}(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^\top \right) p = 0 \\ &\iff \mathbb{E} \left((\mathbf{y} - \bar{\mathbf{y}})^\top p \right)^2 = 0 \end{aligned} \quad (29.145)$$

from which we conclude that the zero-mean scalar random variable $(\mathbf{y} - \bar{y})^\top p$ has zero variance and, hence,

$$(\mathbf{y} - \bar{y})^\top p = 0, \quad \text{in probability} \quad (29.146)$$

This conclusion leads to a contradiction since it implies that

$$\begin{aligned} r_{yx}^\top p &= \mathbb{E}(\mathbf{x} - \bar{x})(\mathbf{y} - \bar{y})^\top p \\ &\stackrel{(29.146)}{=} 0, \quad \text{in probability} \end{aligned} \quad (29.147)$$

which violates the assertion that $r_{yx}^\top p \neq 0$. We conclude that (29.143) holds and the normal equations (29.25) are consistent.

Next we verify that the solution w^o is either unique or there are infinitely many solutions. To begin with, it is clear from (29.25) that the solution is unique whenever R_y is invertible, in which case $w^o = R_y^{-1} r_{yx}$. On the other hand, when R_y is singular, then infinitely many solutions exist. This is because if we let p denote any nontrivial vector in the nullspace of R_y , i.e., $R_y p = 0$, then the vector $w^o + p$ will also satisfy the normal equations (29.25).

The next property we verify is that when infinitely many solutions exist, any solution will continue to lead to the *same* estimator for \mathbf{x} and to the *same* mean-square-error value. Let w^o and w^\bullet denote any two solutions to the normal equations, i.e.,

$$R_y w^o = r_{yx}, \quad R_y w^\bullet = r_{yx} \quad (29.148)$$

The corresponding estimators for \mathbf{x} are denoted by

$$\hat{\mathbf{x}}^o = \bar{x} + (\mathbf{y} - \bar{y})^\top w^o \quad (29.149a)$$

$$\hat{\mathbf{x}}^\bullet = \bar{x} + (\mathbf{y} - \bar{y})^\top w^\bullet \quad (29.149b)$$

Subtracting both equalities in (29.148) gives

$$R_y(w^o - w^\bullet) = 0 \quad (29.150)$$

so that any two solution vectors differ by vectors in the nullspace of R_y , namely,

$$w^\bullet = w^o + p, \quad \text{for some } p \in \mathcal{N}(R_y) \quad (29.151)$$

Moreover, we obtain from (29.150) that

$$\begin{aligned} R_y(w^o - w^\bullet) = 0 &\implies (w^o - w^\bullet)^\top R_y(w^o - w^\bullet) = 0 \\ &\implies (w^o - w^\bullet)^\top \left(\mathbb{E}(\mathbf{y} - \bar{y})(\mathbf{y} - \bar{y})^\top \right) (w^o - w^\bullet) = 0 \\ &\implies \mathbb{E} \left((\mathbf{y} - \bar{y})^\top (w^o - w^\bullet) \right)^2 = 0 \end{aligned} \quad (29.152)$$

which implies that the following zero-mean scalar random variable is equal to zero in probability:

$$\boldsymbol{\alpha} \triangleq (\mathbf{y} - \bar{y})^\top (w^o - w^\bullet) = 0, \quad \text{in probability} \quad (29.153)$$

That is, for any $\epsilon > 0$,

$$\mathbb{P}(|\boldsymbol{\alpha}| \geq \epsilon) = 0 \quad (29.154)$$

Subtracting expressions (29.149a)–(29.149b) gives

$$\begin{aligned} \hat{\mathbf{x}}^o - \hat{\mathbf{x}}^\bullet &= (\mathbf{y} - \bar{y})^\top (w^o - w^\bullet) \\ &\stackrel{(29.153)}{=} 0, \quad \text{in probability} \end{aligned} \quad (29.155)$$

which confirms our claim that different solutions to the normal equations continue to

lead to the same estimator. It is left as an exercise to check that the mean-square-errors corresponding to $\hat{\mathbf{x}}^o$ and $\hat{\mathbf{x}}^\bullet$ also agree with each other — see Prob. 29.4. ■

REFERENCES

- Aitken, A. C. (1935), “On least squares and linear combinations of observations,” *Proceedings of the Royal Society of Edinburgh*, vol. 55, pp. 42–48.
- Cox, D. R. and D. V. Hinkley (1974), *Theoretical Statistics*, Chapman and Hall, NY.
- Gauss, C. F. (1821), *Theoria Combinationis Observationum Erroribus Minimis Obnoxiae*, H. Dieterich, vol. 2, pp. 1–58.
- Kailath, T., A. H. Sayed, and B. Hassibi (2000), *Linear Estimation*, Prentice Hall, NJ.
- Kay, S. (1993), *Fundamentals of Statistical Signal Processing: Estimation Theory*, Prentice Hall, NJ.
- Kendall, M. and A. Stuart (1976–1979), *The Advanced Theory of Statistics*, vols. 1–3, Macmillan, NY.
- Kolmogorov, A. N. (1939), “Sur l’interpolation et extrapolation des suites stationnaires,” *C. R. Acad. Sci.*, vol. 208, p. 2043.
- Kolmogorov, A. N. (1941a), “Stationary sequences in Hilbert space (in Russian),” *Bull. Math. Univ. Moscow*, vol. 2.
- Kolmogorov, A. N. (1941b), “Interpolation and extrapolation of stationary random processes,” *Bull. Acad. Sci. USSR*, vol. 5. [A translation has been published by the RAND Corp., Santa Monica, CA, as Memo. RM-3090-PR, Apr. 1962.]
- Markov, A. A. (1912), *Wahrscheinlichkeitsrechnung*, B. G. Teubner, Leipzig, Berlin.
- Plackett, R. L. (1949), “A historical note on the method of least squares,” *Biometrika*, vol. 36, no. 3–4, pp. 458–4601.
- Plackett, R. L. (1950), “Some theorems in least-squares,” *Biometrika*, vol. 37, no. 1–2, pp. 149–157.
- Rao, C. R. (1973), *Linear Statistical Inference and Its Applications*, Wiley, NY.
- Sayed, A. H. (2003), *Fundamentals of Adaptive Filtering*, Wiley, NJ.
- Sayed, A. H. (2008), *Adaptive Filters*, Wiley, NJ.
- Wiener, N. (1949), *Extrapolation, Interpolation and Smoothing of Stationary Time Series*, Technology Press and Wiley, NY. [Originally published in 1942 as a classified National Defense Research Council Report. Also published under the title *Time Series Analysis* by MIT Press, Cambridge, MA.]
- Wiener, N. and E. Hopf (1931), “On a class of singular integral equations,” *Proc. Prussian Acad. Math. – Phys. Ser.*, p. 696, 1931.
- Wold, H. (1938), *A Study in the Analysis of Stationary Time Series*, Almqvist & Wiksell, Uppsala, Sweden.