

7-4. Mean-square estimation; the orthogonality principle

We shall now be concerned with the question of estimating a r.v. x by a constant or by a function of another r.v. y . This problem will be reexamined in Sec. 8-2 and in Chap. 11. In the following discussion we consider its meaning and introduce the notion of mean square (abbreviation: m.s.) estimation. The formulation and solution of the problem will be in terms of probabilities (conceptual). However, a brief explanation in terms of repeated trials (physical) might be helpful.

Frequency interpretation. A r.v. x is defined on a certain experiment. Its distribution $F(x)$ is given, and so is its value $x(\zeta)$ for each outcome ζ . This does not, of course, mean that if the corresponding physical experiment is performed, one will know in advance the resulting value $x(\zeta)$ of x . The outcome ζ of a particular trial might be any element of \mathcal{S} . The question arises whether, guided by $F(x)$, we could "guess" a value a for $x(\zeta)$. This is the problem of estimating the r.v. x by a constant. Suppose that a is somehow selected. At each trial we commit an error,

$$x(\zeta) - a \quad (7-89)$$

and our problem is to find the particular a that will make this error "small." If by "small" we mean that the average of $x(\zeta) - a$ in a long run of trials should be close to zero,

$$\frac{x(\zeta_1) - a + \cdots + x(\zeta_n) - a}{n} \simeq 0$$

then [see (5-22)] a should equal the expected value of x .

However, depending on the nature of the problem, one might prefer some other criterion for selecting a , for example, the minimization of the average of $|x(\zeta) - a|$. In this case, a should equal the median of x (see Prob. 5-3). In our analysis we shall deal only with m.s. estimations. This means that a should be so selected that the average of

$$[x(\zeta) - a]^2$$

is minimum. This criterion is, in general, useful, but it is primarily chosen because it leads to simple results. We shall soon see that the best a is again the expected value of x ,

$$a = E\{x\} \quad (7-90)$$

The estimation of x can be improved if one has access to the values of another r.v. y . We elaborate: It is assumed that at each trial we "observe" the resulting value $y(\zeta)$ of y and want $x(\zeta)$ estimated on the basis of this observation. If x and y are independent, then knowledge of $y(\zeta)$ is of no help in the estimate of x . In this case x is again estimated by a constant. However, if x and y are not independent, then it might be best to use for an estimate of x not the same number at each trial, but a quantity that depends on the observed $y(\zeta)$. In other words, we want x estimated by a function $g(y)$, and our problem is to find the best $g(y)$.

One might argue that if $y(\zeta)$ is observed, then the particular outcome ζ is known; hence $x(\zeta)$ can be *predicted* exactly. This is not so. The same number $y = y(\zeta)$ might result from several outcomes ζ ,

$$y = y(\zeta_1) = \dots = y(\zeta_n) = \dots \quad (7-91)$$

and for each such ζ the corresponding values of x might be different. Hence, having observed $y(\zeta)$ at a given trial, we cannot, in general, predict $x(\zeta)$, but only estimate† it.

The foregoing reasoning leads to the following conclusion: The fact that $y(\zeta) = y$ is specified means that the outcome ζ of our trial is not any element of \mathcal{S} , but only one of the elements ζ_i in (7-91). In other words, we are asking for an estimate of x in the subset $\{y = y\}$ of our space. In this set, $y(\zeta)$ is a constant, and our problem is to estimate x by the *constant* $g(y(\zeta))$. Changing probabilities into conditional probabilities, we conclude, as in (7-90), that the best m.s. estimate of x is its expected value

$$g(y) = E\{x|y\}$$

We shall soon see that the above loose conclusions can be strictly established in the conceptual world of probabilities.

Mean-square estimation of a random variable by a constant. We start with the following simple but basic problem: Find a constant a such that

$$E\{(x - a)^2\} = \int_{-\infty}^{\infty} (x - a)^2 f(x) dx$$

is minimum. We maintain that

$$a = E\{x\} = \eta_x = \int_{-\infty}^{\infty} xf(x) dx \quad (7-92)$$

Indeed, expanding, we have

$$E\{(x - a)^2\} = a^2 - 2aE\{x\} + E\{x^2\}$$

The derivative with respect to a equals zero for $a = E\{x\}$, and (7-92) follows. Thus the constant η_x has the properties‡ that: The expected value of $x - \eta_x$ equals zero; the expected value of $(x - \eta_x)^2$ is minimum.

Nonlinear mean-square estimation of y in terms of x . We now want to estimate the r.v. y by a suitable function $g(x)$ of x so that the m.s. estimation error

$$E\{[y - g(x)]^2\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} [y - g(x)]^2 f(x, y) dx dy \quad (7-93)$$

is minimum (we reversed the role of x and y).

† In the literature the expression "prediction of x in terms of y " is often used; from the above we see that "estimation" is a more appropriate term.

‡ The above result corresponds to the well-known fact that the moment of inertia with respect to the center of gravity is smaller than with respect to any other point.

Theorem. The function $g(x)$ that minimizes (7-93) is the conditional expected value of y , assuming x :

$$g(x) = E\{y|x\} \quad (7-94)$$

Proof. Since

$$f(x, y) = f(y|x)f(x)$$

we have

$$E\{[y - g(x)]^2\} = \int_{-\infty}^{\infty} f(x) \int_{-\infty}^{\infty} [y - g(x)]^2 f(y|x) dy dx$$

The integrand above is nonnegative; therefore, in order to minimize the double integral, it suffices to minimize

$$\int_{-\infty}^{\infty} [y - g(x)]^2 f(y|x) dy$$

for every x . For a given x , this integral is the second moment of the conditional density $f(y|x)$ with respect to the constant $g(x)$. As we know from (7-92), this moment is minimum if

$$g(x) = \int_{-\infty}^{\infty} y f(y|x) dy = E\{y|x\}$$

and (7-94) follows.

Thus (7-94) is a simple extension of (7-92) in the probability space conditioned by $\{x = x\}$. This conclusion can also be drawn, after some thought, from [see (7-59)]

$$E\{[y - g(x)]^2\} = E\{E\{[y - g(x)]^2|x\}\}$$

The function

$$g(x) = E\{y|x\}$$

is known as regression curve (Fig. 7-18). It is the locus of the centers of gravity of the masses on the strips $(x, x + dx)$. If these masses are near $g(x)$, then the m.s. error

$$E\{[y - E\{y|x\}]^2\} \quad (7-95)$$

is small.

Independent Random Variables. If x and y are independent, then (see page 182)

$$E\{y|x\} = E\{y\}$$

Hence the best m.s. estimate of y in terms of x is $E\{y\}$. Thus knowledge of x does not help in the estimation of y .

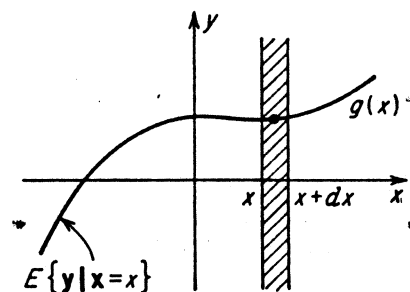


Fig. 7-18

Linear Mean-square Estimation; the Orthogonality Principle

The solution $E\{y|x\}$ of the nonlinear estimation problem looks simple enough. However, the actual evaluation of $E\{y|x\}$ is not simple at all. One must determine this function for every x . The difficulties become severe if more than one r.v. are involved (Chap. 11). A much easier problem is the estimation of y by a linear function

$$ax + b$$

of x . We now seek not a function, but merely the two constants a and b that minimize

$$E\{[y - (ax + b)]^2\}$$

The resulting error is of course larger than the corresponding error in the nonlinear estimation; however, this is often compensated by the simplicity of the solution.

Theorem. The constants a and b that minimize the m.s. error

$$e = E\{[y - (ax + b)]^2\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (y - ax - b)^2 f(x, y) dx dy \quad (7-96)$$

are given by

$$a = \frac{r\sigma_y}{\sigma_x} \quad b = E\{y\} - aE\{x\} \quad (7-97)$$

and the resulting minimum error e_m by

$$e_m = \sigma_y^2(1 - r^2) \quad (7-98)$$

where r is the correlation coefficient of x and y [see (7-66)].

Proof. Suppose that a is specified. The value of b that minimizes e is the best m.s. estimate of the r.v. $y - ax$ by a constant; hence [see (7-92)]

$$b = E\{y - ax\} = \eta_y - a\eta_x$$

With b so determined, we now have

$$E\{(y - ax - b)^2\} = E\{[(y - \eta_y) - a(x - \eta_x)]^2\} = \sigma_y^2 - 2r\sigma_x\sigma_y a + \sigma_x^2 a^2$$

The last quantity is minimum for

$$a = \frac{r\sigma_x\sigma_y}{\sigma_x^2} = \frac{r\sigma_y}{\sigma_x}$$

and (7-97) follows. Inserting the value of a in the above quadratic, we find

$$e_m = \sigma_y^2 - 2r^2\sigma_y^2 + r^2\sigma_y^2 = \sigma_y^2(1 - r^2)$$

$$r = \frac{E\{(x - \eta_x)(y - \eta_y)\}}{\sqrt{E\{(x - \eta_x)^2\} E\{(y - \eta_y)^2\}}} = \frac{\mu_{xy}}{\sigma_x \sigma_y}$$

The above result will now be stated in a more basic form. We shall assume for simplicity that

$$E\{x\} = E\{y\} = 0$$

Orthogonality principle. The constant a that minimizes the m.s. error

$$e = E\{(y - ax)^2\}$$

is such that $y - ax$ is orthogonal to x ; that is,

$$E\{(y - ax)x\} = 0 \quad (7-99)$$

and the minimum m.s. error is given by

$$e_m = E\{(y - ax)y\} \quad (7-100)$$

Proof. The above can be deduced from (7-97) and (7-98); it will be instructive, however, to give a second proof. This proof can be simply extended to complex r.v. (see Sec. 11-2), whereas the differentiation presents certain complications.

Suppose that a is such that $E\{(y - ax)x\} = 0$. We maintain that the resulting error e is minimum. Indeed, for any A , we have

$$\begin{aligned} E\{(y - Ax)^2\} &= E\{[(y - ax) + (a - A)x]^2\} \\ &= E\{(y - ax)^2\} + 2(a - A)E\{(y - ax)x\} + (a - A)^2E\{x^2\} \end{aligned}$$

But the second term in the last expression is zero, and the last non-negative; hence

$$E\{(y - Ax)^2\} \geq E\{(y - ax)^2\}$$

and our statement is proved. The minimum error is given by

$$e_m = E\{(y - ax)^2\} = E\{(y - ax)y\} - aE\{(y - ax)x\}$$

and (7-100) follows because the last term is zero.

From (7-99) we conclude that

$$a = \frac{E\{xy\}}{E\{x^2\}} \quad (7-101)$$

Inserting into (7-100), we obtain

$$e_m = E\{y^2\} - aE\{xy\} = E\{y^2\} - \frac{E^2\{xy\}}{E\{x^2\}} \quad (7-102)$$

The m.s. error can also be written in the form

$$e_m = E\{y^2\} - E\{(ax)^2\} \quad (7-103)$$

We remark that

$$e_m \geq E\{[y - E\{y|x\}]^2\}$$

We shall now extend briefly the results of Sec. 7-4 to several r.v. This discussion will be resumed in Chap. 11. We are given the $n + 1$ r.v.

$$x_0, x_1, \dots, x_n \quad (8-29)$$

and we want to estimate x_0 by a function $g(x_1, \dots, x_n)$ of the other r.v. so as to minimize the m.s. error:

$$E\{[x_0 - g(x_1, \dots, x_n)]^2\} \quad (8-30)$$

Reasoning as in (7-94), we can easily show that

$$g(x_1, \dots, x_n) = E\{x_0|x_1, \dots, x_n\} \quad (8-31)$$

The above expected value is given by (8-17). Thus, to solve the non-linear m.s. estimation problem, we need to know the joint density of the r.v. x_0, \dots, x_n .

Linear mean-square estimation. The estimation problem is considerably simplified if one seeks an estimate of x_0 by a linear combination of x_1, \dots, x_n . In this case the problem is to find n constants a_1, \dots, a_n such that the m.s. error

$$e = E\{[x_0 - (a_1x_1 + \dots + a_nx_n)]^2\} \quad (8-32)$$

is minimum. It turns out that these constants can be determined in terms of the second moments

$$R_{ij} = E\{x_ix_j\}$$

of the given r.v. If $E\{x_i\} = 0$, then R_{ij} is the covariance of the r.v. x_i and x_j .

Orthogonality Principle. The constants a_i that minimize e are such that the error

$$x_0 - (a_1x_1 + \dots + a_nx_n)$$

is orthogonal to x_1, \dots, x_n ; that is,

$$E\{[x_0 - (a_1x_1 + \dots + a_nx_n)]x_i\} = 0 \quad i = 1, \dots, n \quad (8-33)$$

Proof. The m.s. error e is a function of a_1, \dots, a_n , and to minimize it we differentiate with respect to a_i :

$$\frac{\partial e(a_1, \dots, a_n)}{\partial a_i} = \frac{\partial E\{[x_0 - (a_1x_1 + \dots + a_nx_n)]^2\}}{\partial a_i} = 0 \quad i = 1, \dots, n$$

Writing the above expected value as an integral of the form (8-15), we see that the order of differentiation and expected value can be interchanged; the result is

$$\frac{\partial e}{\partial a_i} = -2E\{[x_0 - (a_1x_1 + \dots + a_nx_n)]x_i\} = 0$$

and (8-33) follows.

It is easy to see by expanding the square in (8-32) and using (8-33) that the minimum m.s. error is given by

$$\begin{aligned} e_m &= E\{[x_0 - (a_1x_1 + \dots + a_nx_n)]x_0\} \\ &= R_{00} - (a_1R_{01} + \dots + a_nR_{0n}) \end{aligned} \quad (8-34)$$