

# 16. Mean Square Estimation

Given some information that is related to an unknown quantity of interest, the problem is to obtain a good estimate for the unknown in terms of the observed data.

Suppose  $X_1, X_2, \dots, X_n$  represent a sequence of random variables about whom one set of observations are available, and  $Y$  represents an unknown random variable. The problem is to obtain a good estimate for  $Y$  in terms of the observations  $X_1, X_2, \dots, X_n$ .

Let

$$\hat{Y} = \varphi(X_1, X_2, \dots, X_n) = \varphi(\underline{X}) \quad (16-1)$$

represent such an estimate for  $Y$ .

Note that  $\varphi(\cdot)$  can be a linear or a nonlinear function of the observation  $X_1, X_2, \dots, X_n$ . Clearly

$$\varepsilon(\underline{X}) = Y - \hat{Y} = Y - \varphi(\underline{X}) \quad (16-2)$$

represents the error in the above estimate, and  $|\varepsilon|^2$  the square of

the error. Since  $\varepsilon$  is a random variable,  $E\{|\varepsilon|^2\}$  represents the mean square error. One strategy to obtain a good estimator would be to minimize the mean square error by varying over all possible forms of  $\varphi(\cdot)$ , and this procedure gives rise to the Minimization of the Mean Square Error (MMSE) criterion for estimation. Thus under MMSE criterion, the estimator  $\varphi(\cdot)$  is chosen such that the mean square error  $E\{|\varepsilon|^2\}$  is at its minimum.

Next we show that the conditional mean of  $Y$  given  $\underline{X}$  is the best estimator in the above sense.

**Theorem 1:** Under MMSE criterion, the best estimator for the unknown  $Y$  in terms of  $X_1, X_2, \dots, X_n$  is given by the conditional mean of  $Y$  given  $\underline{X}$ . Thus

$$\hat{Y} = \varphi(\underline{X}) = E\{Y | \underline{X}\}. \quad (16-3)$$

**Proof:** Let  $\hat{Y} = \varphi(\underline{X})$  represent an estimate of  $Y$  in terms of

$\underline{X} = (X_1, X_2, \dots, X_n)$ . Then the error  $\varepsilon = Y - \hat{Y}$ , and the mean square error is given by

$$\sigma_\varepsilon^2 = E\{|\varepsilon|^2\} = E\{|Y - \hat{Y}|^2\} = E\{|Y - \varphi(\underline{X})|^2\} \quad (16-4) \quad \text{PILLAI}^2$$

Since

$$E[z] = E_X[E_z\{z | \underline{X}\}] \quad (16-5)$$

we can rewrite (16-4) as

$$\sigma_\varepsilon^2 = E\left\{\underbrace{|Y - \varphi(\underline{X})|^2}_z\right\} = E_X\left[E_Y\left\{\underbrace{|Y - \varphi(\underline{X})|^2}_z \mid \underline{X}\right\}\right]$$

where the inner expectation is with respect to  $Y$ , and the outer one is with respect to  $\underline{X}$ .

Thus

$$\begin{aligned} \sigma_\varepsilon^2 &= E[E\{|Y - \varphi(\underline{X})|^2 \mid \underline{X}\}] \\ &= \int_{-\infty}^{+\infty} E\{|Y - \varphi(\underline{X})|^2 \mid \underline{X}\} f_X(\underline{X}) dx. \end{aligned} \quad (16-6)$$

To obtain the best estimator  $\varphi$ , we need to minimize  $\sigma_\varepsilon^2$  in (16-6) with respect to  $\varphi$ . In (16-6), since  $f_X(\underline{X}) \geq 0$ ,  $E\{|Y - \varphi(\underline{X})|^2 \mid \underline{X}\} \geq 0$ , and the variable  $\varphi$  appears only in the integrand term, minimization of the mean square error  $\sigma_\varepsilon^2$  in (16-6) with respect to  $\varphi$  is equivalent to minimization of  $E\{|Y - \varphi(\underline{X})|^2 \mid \underline{X}\}$  with respect to  $\varphi$ .

Since  $\underline{X}$  is fixed at some value,  $\varphi(\underline{X})$  is no longer random, and hence minimization of  $E\{|Y - \varphi(\underline{X})|^2 | \underline{X}\}$  is equivalent to

$$\frac{\partial}{\partial \varphi} E\{|Y - \varphi(\underline{X})|^2 | \underline{X}\} = 0. \quad (16-7)$$

This gives

$$E\{Y - \varphi(\underline{X}) | \underline{X}\} = 0$$

or

$$E\{Y | \underline{X}\} - E\{\varphi(\underline{X}) | \underline{X}\} = 0. \quad (16-8)$$

But

$$E\{\varphi(\underline{X}) | \underline{X}\} = \varphi(\underline{X}), \quad (16-9)$$

since when  $\underline{X} = \underline{x}$ ,  $\varphi(\underline{X})$  is a fixed number  $\varphi(\underline{x})$ . Using (16-9)

in (16-8) we get the desired estimator to be

$$\hat{Y} = \varphi(\underline{X}) = E\{Y | \underline{X}\} = E\{Y | X_1, X_2, \dots, X_n\}. \quad (16-10)$$

Thus the conditional mean of  $Y$  given  $X_1, X_2, \dots, X_n$  represents the best estimator for  $Y$  that minimizes the mean square error.

The minimum value of the mean square error is given by

$$\begin{aligned} \sigma_{\min}^2 &= E\{|Y - E(Y | X)|^2\} = E[\underbrace{E\{|Y - E(Y | X)|^2 | \underline{X}\}}_{\text{var}(Y|\underline{X})}] \\ &= E\{\text{var}(Y | \underline{X})\} \geq 0. \end{aligned} \quad (16-11)$$

As an example, suppose  $Y = X^3$  is the unknown. Then the best MMSE estimator is given by

$$\hat{Y} = E\{Y | X\} = E\{X^3 | X\} = X^3. \quad (16-12)$$

Clearly if  $Y = X^3$ , then indeed  $\hat{Y} = X^3$  is the best estimator for  $Y$

in terms of  $X$ . Thus the best estimator can be nonlinear.

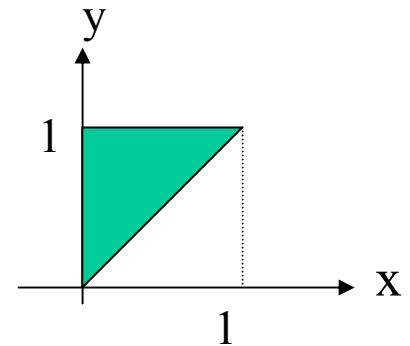
Next, we will consider a less trivial example.

**Example :** Let

$$f_{x,y}(x,y) = \begin{cases} kxy, & 0 < x < y < 1 \\ 0 & \text{otherwise,} \end{cases}$$

where  $k > 0$  is a suitable normalization constant. To determine the best estimate for  $Y$  in terms of  $X$ , we need  $f_{y|x}(y|x)$ .

$$\begin{aligned} f_x(x) &= \int_x^1 f_{x,y}(x,y) dy = \int_x^1 kxy dy \\ &= \frac{kxy^2}{2} \Big|_x^1 = \frac{kx(1-x^2)}{2}, \quad 0 < x < 1. \end{aligned}$$



Thus

$$f_{y|x}(y|x) = \frac{f_{x,y}(x,y)}{f_x(x)} = \frac{kxy}{kx(1-x^2)/2} = \frac{2y}{1-x^2}; \quad 0 < x < y < 1. \quad (16-13)$$

Hence the best MMSE estimator is given by

$$\begin{aligned}
\hat{Y} = \varphi(X) &= E\{Y | \underline{X}\} = \int_x^1 y f_{y|x}(y | x) dy \\
&= \int_x^1 y \frac{2y}{1-x^2} dy = \frac{2}{1-x^2} \int_x^1 y^2 dy \\
&= \frac{2}{3} \frac{y^3}{1-x^2} \Big|_x^1 = \frac{2}{3} \frac{1-x^3}{1-x^2} = \frac{2}{3} \frac{(1+x+x^2)}{1-x^2}.
\end{aligned} \tag{16-14}$$

Once again the best estimator is nonlinear. In general the best estimator  $E\{Y | \underline{X}\}$  is difficult to evaluate, and hence next we will examine the special subclass of best linear estimators.

## Best Linear Estimator

In this case the estimator  $\hat{y}$  is a linear function of the observations  $X_1, X_2, \dots, X_n$ . Thus

$$\hat{Y}_l = a_1 X_1 + a_2 X_2 + \dots + a_n X_n = \sum_{i=1}^n a_i X_i. \tag{16-15}$$

where  $a_1, a_2, \dots, a_n$  are unknown quantities to be determined. The mean square error is given by  $(\varepsilon = Y - \hat{Y}_l)$

$$E\{|\varepsilon|^2\} = E\{|Y - \hat{Y}_l|^2\} = E\{|Y - \sum a_i X_i|^2\} \quad (16-16)$$

and under the MMSE criterion  $a_1, a_2, \dots, a_n$  should be chosen so that the mean square error  $E\{|\varepsilon|^2\}$  is at its minimum possible value. Let  $\sigma_n^2$  represent that minimum possible value. Then

$$\sigma_n^2 = \min_{a_1, a_2, \dots, a_n} E\{|Y - \sum_{i=1}^n a_i X_i|^2\}. \quad (16-17)$$

To minimize (16-16), we can equate

$$\frac{\partial}{\partial a_k} E\{|\varepsilon|^2\} = 0, \quad k = 1, 2, \dots, n. \quad (16-18)$$

This gives

$$\frac{\partial}{\partial a_k} E\{|\varepsilon|^2\} = E\left\{\frac{\partial |\varepsilon|^2}{\partial a_k}\right\} = 2E\left[\varepsilon \left\{\frac{\partial \varepsilon}{\partial a_k}\right\}^*\right] = 0. \quad (16-19)$$

But



$$\frac{\partial \varepsilon}{\partial a_k} = \frac{\partial(Y - \sum_{i=1}^n a_i X_i)}{\partial a_k} = \frac{\partial Y}{\partial a_k} - \frac{\partial(\sum_{i=1}^n a_i X_i)}{\partial a_k} = -X_k. \quad (16-20)$$

Substituting (16-19) in to (16-18), we get

$$\frac{\partial E\{|\varepsilon|^2\}}{\partial a_k} = -2E\{\varepsilon X_k^*\} = 0,$$

or the best linear estimator must satisfy

$$E\{\varepsilon X_k^*\} = 0, \quad k = 1, 2, \dots, n. \quad (16-21)$$

Notice that in (16-21),  $\varepsilon$  represents the estimation error  $(Y - \sum_{i=1}^n a_i X_i)$ , and  $X_k$ ,  $k = 1 \rightarrow n$  represents the data. Thus from (16-21), the error  $\varepsilon$  is orthogonal to the data  $X_k$ ,  $k = 1 \rightarrow n$  for the best linear estimator. This is the **orthogonality principle**.

In other words, in the linear estimator (16-15), the unknown constants  $a_1, a_2, \dots, a_n$  must be selected such that the error

$\varepsilon = Y - \sum_{i=1}^n a_i X_i$  is orthogonal to every data  $X_1, X_2, \dots, X_n$  for the best linear estimator that minimizes the mean square error.

Interestingly a general form of the orthogonality principle holds good in the case of nonlinear estimators also.

**Nonlinear Orthogonality Rule:** Let  $h(\underline{X})$  represent *any* functional form of the data and  $E\{Y | \underline{X}\}$  the best estimator for  $Y$  given  $\underline{X}$ . With  $e = Y - E\{Y | \underline{X}\}$  we shall show that

$$E\{eh(\underline{X})\} = 0, \quad (16-22)$$

implying that

$$e = Y - E\{Y | \underline{X}\} \perp h(\underline{X}).$$

This follows since

$$\begin{aligned} E\{eh(\underline{X})\} &= E\{(Y - E[Y | \underline{X}])h(\underline{X})\} \\ &= E\{Yh(\underline{X})\} - E\{E[Y | \underline{X}]h(\underline{X})\} \\ &= E\{Yh(\underline{X})\} - E\{E[Yh(\underline{X}) | \underline{X}]\} \\ &= E\{Yh(\underline{X})\} - E\{Yh(\underline{X})\} = 0. \end{aligned}$$

Thus in the nonlinear version of the orthogonality rule the error is orthogonal to *any* functional form of the data.

The orthogonality principle in (16-20) can be used to obtain the unknowns  $a_1, a_2, \dots, a_n$  in the linear case.

For example suppose  $n = 2$ , and we need to estimate  $Y$  in terms of  $X_1$  and  $X_2$  linearly. Thus

$$\hat{Y}_l = a_1 X_1 + a_2 X_2$$

From (16-20), the orthogonality rule gives

$$E\{\varepsilon X_1^*\} = E\{(Y - a_1 X_1 - a_2 X_2) X_1^*\} = 0$$

$$E\{\varepsilon X_2^*\} = E\{(Y - a_1 X_1 - a_2 X_2) X_2^*\} = 0$$

Thus

$$E\{|X_1|^2\}a_1 + E\{X_2 X_1^*\}a_2 = E\{Y X_1^*\}$$

$$E\{X_1 X_2^*\}a_1 + E\{|X_2|^2\}a_2 = E\{Y X_2^*\}$$

or

$$\begin{pmatrix} E\{|X_1|^2\} & E\{X_2 X_1^*\} \\ E\{X_1 X_2^*\} & E\{|X_2|^2\} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} E\{Y X_1^*\} \\ E\{Y X_2^*\} \end{pmatrix} \quad (16-23)$$

(16-23) can be solved to obtain  $a_1$  and  $a_2$  in terms of the cross-correlations.

The minimum value of the mean square error  $\sigma_n^2$  in (16-17) is given by

$$\begin{aligned} \sigma_n^2 &= \min_{a_1, a_2, \dots, a_n} E\{|\varepsilon|^2\} \\ &= \min_{a_1, a_2, \dots, a_n} E\{\varepsilon \varepsilon^*\} = \min_{a_1, a_2, \dots, a_n} E\{\varepsilon (Y - \sum_{i=1}^n a_i X_i)^*\} \\ &= \min_{a_1, a_2, \dots, a_n} E\{\varepsilon Y^*\} - \min_{a_1, a_2, \dots, a_n} \sum_{i=1}^n a_i E\{\varepsilon X_i^*\}. \end{aligned} \quad (16-24)$$

But using (16-21), the second term in (16-24) is zero, since the error is orthogonal to the data  $X_i$ , where  $a_1, a_2, \dots, a_n$  are chosen to be optimum. Thus the minimum value of the mean square error is given by

$$\begin{aligned}\sigma_n^2 &= E\{\varepsilon Y^*\} = E\left\{\left(Y - \sum_{i=1}^n a_i X_i\right)Y^*\right\} \\ &= E\{|Y|^2\} - \sum_{i=1}^n a_i E\{X_i Y^*\}\end{aligned}\quad (16-25)$$

where  $a_1, a_2, \dots, a_n$  are the optimum values from (16-21).

Since the linear estimate in (16-15) is only a special case of the general estimator  $\varphi(\underline{X})$  in (16-1), the best linear estimator that satisfies (16-20) cannot be superior to the best nonlinear estimator  $E\{Y | \underline{X}\}$ . Often the best linear estimator will be inferior to the best estimator in (16-3).

This raises the following question. Are there situations in which the best estimator in (16-3) also turns out to be linear? In those situations it is enough to use (16-21) and obtain the best linear estimators, since they also represent the best global estimators. Such is the case if  $Y$  and  $X_1, X_2, \dots, X_n$  are distributed as jointly Gaussian.

We summarize this in the next theorem and prove that result.

**Theorem2:** If  $X_1, X_2, \dots, X_n$  and  $Y$  are jointly Gaussian zero

mean random variables, then the best estimate for  $Y$  in terms of  $X_1, X_2, \dots, X_n$  is always linear.

**Proof :** Let

$$\hat{Y} = \varphi(X_1, X_2, \dots, X_n) = E\{Y | \underline{X}\} \quad (16-26)$$

represent the best (possibly nonlinear) estimate of  $Y$ , and

$$\hat{Y}_l = \sum_{i=1}^n a_i X_i \quad (16-27)$$

the best linear estimate of  $Y$ . Then from (16-21)

$$\varepsilon \triangleq Y - Y_l = Y - \sum_{i=1}^n a_i X_i \quad (16-28)$$

is orthogonal to the data  $X_k$ ,  $k = 1 \rightarrow n$ . Thus

$$E\{\varepsilon X_k^*\} = 0, \quad k = 1 \rightarrow n. \quad (16-29)$$

Also from (16-28),

$$E\{\varepsilon\} = E\{Y\} - \sum_{i=1}^n a_i E\{X_i\} = 0. \quad (16-30)$$

Using (16-29)-(16-30), we get

$$E\{\varepsilon X_k^*\} = E\{\varepsilon\}E\{X_k^*\} = 0, \quad k = 1 \rightarrow n. \quad (16-31)$$

From (16-31), we obtain that  $\varepsilon$  and  $X_k$  are zero mean uncorrelated random variables for  $k = 1 \rightarrow n$ . But  $\varepsilon$  itself represents a Gaussian random variable, since from (16-28) it represents a linear combination of a set of jointly Gaussian random variables. Thus  $\varepsilon$  and  $\underline{X}$  are jointly Gaussian and uncorrelated random variables. As a result,  $\varepsilon$  and  $\underline{X}$  are independent random variables. Thus from their independence

$$E\{\varepsilon \mid \underline{X}\} = E\{\varepsilon\}. \quad (16-32)$$

But from (16-30),  $E\{\varepsilon\} = 0$ , and hence from (16-32)

$$E\{\varepsilon \mid \underline{X}\} = 0. \quad (16-33)$$

Substituting (16-28) into (16-33), we get

$$E\{\varepsilon \mid \underline{X}\} = E\left\{Y - \sum_{i=1}^n a_i X_i \mid \underline{X}\right\} = 0$$

or

$$E\{Y | \underline{X}\} = E\left\{\sum_{i=1}^n a_i X_i | \underline{X}\right\} = \sum_{i=1}^n a_i X_i = Y_l. \quad (16-34)$$

From (16-26),  $E\{Y | \underline{X}\} = \varphi(\underline{x})$  represents the best possible estimator, and from (16-28),  $\sum_{i=1}^n a_i X_i$  represents the best linear estimator. Thus the best linear estimator is also the best possible overall estimator in the Gaussian case.

Next we turn our attention to prediction problems using linear estimators.

## Linear Prediction

Suppose  $X_1, X_2, \dots, X_n$  are known and  $X_{n+1}$  is unknown. Thus  $Y = X_{n+1}$ , and this represents a one-step prediction problem. If the unknown is  $X_{n+k}$ , then it represents a  $k$ -step ahead prediction problem. Returning back to the one-step predictor, let  $\hat{X}_{n+1}$  represent the best linear predictor. Then



$$\hat{X}_{n+1} \triangleq -\sum_{i=1}^n a_i X_i, \quad (16-35)$$

where the error

$$\begin{aligned} \varepsilon_n &= X_{n+1} - \hat{X}_{n+1} = X_{n+1} + \sum_{i=1}^n a_i X_i \\ &= a_1 X_1 + a_2 X_2 + \cdots + a_n X_n + X_{n+1} \\ &= \sum_{i=1}^{n+1} a_i X_i, \quad a_{n+1} = 1, \end{aligned} \quad (16-36)$$

is orthogonal to the data, i.e.,

$$E\{\varepsilon_n X_k^*\} = 0, \quad k = 1 \rightarrow n. \quad (16-37)$$

Using (16-36) in (16-37), we get

$$E\{\varepsilon_n X_k^*\} = \sum_{i=1}^{n+1} a_i E\{X_i X_k^*\} = 0, \quad k = 1 \rightarrow n. \quad (16-38)$$

Suppose  $X_i$  represents the sample of a wide sense stationary 17

stochastic process  $X(t)$  so that

$$E\{X_i X_k^*\} = R(i-k) = r_{i-k} = r_{k-i}^* \quad (16-39)$$

Thus (16-38) becomes

$$E\{\varepsilon_n X_k^*\} = \sum_{i=1}^{n+1} a_i r_{i-k} = 0, \quad a_{n+1} = 1, \quad k = 1 \rightarrow n. \quad (16-40)$$

Expanding (16-40) for  $k = 1, 2, \dots, n$ , we get the following set of linear equations.

$$\begin{aligned} a_1 r_0 + a_2 r_1 + a_3 r_2 + \dots + a_n r_{n-1} + r_n &= 0 \leftarrow k = 1 \\ a_1 r_1^* + a_2 r_0 + a_3 r_1 + \dots + a_n r_{n-2} + r_{n-1} &= 0 \leftarrow k = 2 \\ \vdots & \\ a_1 r_{n-1}^* + a_2 r_{n-2}^* + a_3 r_{n-3}^* + \dots + a_n r_0 + r_1 &= 0 \leftarrow k = n. \end{aligned} \quad (16-41)$$

Similarly using (16-25), the minimum mean square error is given by

$$\begin{aligned}
\sigma_n^2 &= E\{|\varepsilon|^2\} = E\{\varepsilon_n Y^*\} = E\{\varepsilon_n X_{n+1}^*\} \\
&= E\left\{\left(\sum_{i=1}^{n+1} a_i X_i\right) X_{n+1}^*\right\} = \sum_{i=1}^{n+1} a_i r_{n+1-i}^* \\
&= a_1 r_n^* + a_2 r_{n-1}^* + a_3 r_{n-2}^* + \cdots + a_n r_1 + r_0.
\end{aligned} \tag{16-42}$$

The  $n$  equations in (16-41) together with (16-42) can be represented as

$$\begin{pmatrix} r_0 & r_1 & r_2 & \cdots & r_n \\ r_1^* & r_0 & r_1 & \cdots & r_{n-1} \\ r_2^* & r_1^* & r_0 & \cdots & r_{n-2} \\ & & \vdots & & \\ r_{n-1}^* & r_{n-2}^* & \cdots & r_0 & r_1 \\ r_n^* & r_{n-1}^* & \cdots & r_1^* & r_0 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ \sigma_n^2 \end{pmatrix}. \tag{16-43}$$

Let

$$T_n = \begin{pmatrix} r_0 & r_1 & r_2 & \cdots & r_n \\ r_1^* & r_0 & r_1 & \cdots & r_{n-1} \\ & & \vdots & & \\ r_n^* & r_{n-1}^* & \cdots & r_1^* & r_0 \end{pmatrix}. \quad (16-44)$$

Notice that  $T_n$  is Hermitian Toeplitz and positive definite. Using (16-44), the unknowns in (16-43) can be represented as

$$\begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \\ 1 \end{pmatrix} = T_n^{-1} \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ \sigma_n^2 \end{pmatrix} = \sigma_n^2 \begin{pmatrix} \text{Last} \\ \text{column} \\ \text{of} \\ T_n^{-1} \end{pmatrix} \quad (16-45)$$

Let

$$T_n^{-1} = \begin{pmatrix} T_n^{11} & T_n^{12} & \cdots & T_n^{1,n+1} \\ T_n^{21} & T_n^{22} & \cdots & T_n^{2,n+1} \\ & & \vdots & \\ T_n^{n+1,1} & T_n^{n+1,2} & \cdots & T_n^{n+1,n+1} \end{pmatrix}. \quad (16-46)$$

Then from (16-45),

$$\begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \\ 1 \end{pmatrix} = \sigma_n^2 \begin{pmatrix} T_n^{1,n+1} \\ T_n^{2,n+1} \\ \vdots \\ T_n^{n+1,n+1} \end{pmatrix}. \quad (16-47)$$

Thus

$$\sigma_n^2 = \frac{1}{T_n^{n+1,n+1}} > 0, \quad (16-48)$$

and

$$\begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} = \frac{1}{T_n^{n+1,n+1}} \begin{pmatrix} T_n^{1,n+1} \\ T_n^{2,n+1} \\ \vdots \\ T_n^{n+1,n+1} \end{pmatrix}. \quad (16-49)$$

Eq. (16-49) represents the best linear predictor coefficients, and they can be evaluated from the last column of  $T_n$  in (16-45). Using these, The best one-step ahead predictor in (16-35) taken the form

$$\hat{X}_{n+1} = -\left(\frac{1}{T_n^{n+1,n+1}}\right) \sum_{i=1}^n (T_n^{i,n+1}) X_i. \quad (16-50)$$

and from (16-48), the minimum mean square error is given by the  $(n+1, n+1)$  entry of  $T_n^{-1}$ .

From (16-36), since the one-step linear prediction error

$$\varepsilon_n = X_{n+1} + a_n X_n + a_{n-1} X_{n-1} + \cdots + a_1 X_1, \quad (16-51)$$

we can represent (16-51) formally as follows

$$X_{n+1} \rightarrow \boxed{1 + a_n z^{-1} + a_{n-1} z^{-2} + \cdots + a_1 z^{-n}} \rightarrow \varepsilon_n$$

Thus, let

$$A_n(z) = 1 + a_n z^{-1} + a_{n-1} z^{-2} + \cdots + a_1 z^{-n}, \quad (16-52)$$

them from the above figure, we also have the representation

$$\varepsilon_n \rightarrow \boxed{\frac{1}{A_n(z)}} \rightarrow X_{n+1}.$$

The filter

$$H(z) = \frac{1}{A_n(z)} = \frac{1}{1 + a_n z^{-1} + a_{n-1} z^{-2} + \cdots + a_1 z^{-n}} \quad (16-53)$$

represents an  $AR(n)$  filter, and this shows that linear prediction leads to an auto regressive ( $AR$ ) model.

The polynomial  $A_n(z)$  in (16-52)-(16-53) can be simplified using (16-43)-(16-44). To see this, we rewrite  $A_n(z)$  as

$$A_n(z) = a_1 z^{-n} + a_2 z^{-(n-1)} + \cdots + a_{n-1} z^{-2} + a_n z^{-1} + 1$$

$$= [z^{-n}, z^{-(n-1)}, \cdots, z^{-1}, 1] \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \\ 1 \end{bmatrix} = [z^{-n}, z^{-(n-1)}, \cdots, z^{-1}, 1] T_n^{-1} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \sigma_n^2 \end{bmatrix} \quad (16-54)$$

To simplify (16-54), we can make use of the following matrix identity

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} I & -AB \\ 0 & I \end{bmatrix} = \begin{bmatrix} A & 0 \\ C & D - CA^{-1}B \end{bmatrix}. \quad (16-55)$$



Taking determinants, we get

$$\begin{vmatrix} A & B \\ C & D \end{vmatrix} = |A| |D - CA^{-1}B|. \quad (16-56)$$

In particular if  $D \equiv 0$ , we get

$$|CA^{-1}B| = \frac{(-1)^n}{|A|} \begin{vmatrix} A & B \\ C & 0 \end{vmatrix}. \quad (16-57)$$

Using (16-57) in (16-54), with

$$C = [z^{-n}, z^{-(n-1)}, \dots, z^{-1}, 1], \quad A = T_n, \quad B = \begin{bmatrix} 0 \\ \vdots \\ \sigma_n^2 \end{bmatrix}$$

we get

$$A_n(z) = \frac{(-1)^n}{|T_n|} \left| \begin{array}{cccc|c} & & & & 0 \\ & & & & 0 \\ & & & & \vdots \\ & & & & 0 \\ & & & & \sigma_n^2 \\ \hline z^{-n}, \dots, z^{-1}, 1 & & & & 0 \end{array} \right| = \frac{\sigma_n^2}{|T_n|} \left| \begin{array}{ccccc} r_0 & r_1 & r_2 & \cdots & r_n \\ r_1^* & r_0 & r_1 & \cdots & r_{n-1} \\ & & \vdots & & \\ r_{n-1}^* & r_{n-2}^* & \cdots & r_0 & r_1 \\ z^{-n} & z^{-(n-1)} & \cdots & z^{-1} & 1 \end{array} \right|. \quad (16-58)$$

Referring back to (16-43), using Cramer's rule to solve for  $a_{n+1}$  ( $=1$ ), we get

$$a_{n+1} = \frac{\sigma_n^2 \left| \begin{array}{ccc} r_0 & \cdots & r_{n-1} \\ & \vdots & \\ r_{n-1} & \cdots & r_0 \end{array} \right|}{|T_n|} = \sigma_n^2 \frac{|T_{n-1}|}{|T_n|} = 1$$

or

$$\sigma_n^2 = \frac{|T_n|}{|T_{n-1}|} > 0. \quad (16-59)$$

Thus the polynomial (16-58) reduces to

$$A_n(z) = \frac{1}{|T_{n-1}|} \begin{vmatrix} r_0 & r_1 & r_2 & \cdots & r_n \\ r_1^* & r_0 & r_1 & \cdots & r_{n-1} \\ & & \vdots & & \\ r_{n-1}^* & r_{n-2}^* & \cdots & r_0 & r_1 \\ z^{-n} & z^{-(n-1)} & \cdots & z^{-1} & 1 \end{vmatrix} \quad (16-60)$$

$$= 1 + a_n z^{-1} + a_{n-1} z^{-2} + \cdots + a_1 z^{-n}.$$

The polynomial  $A_n(z)$  in (16-53) can be alternatively represented as

in (16-60), and  $H(z) = \frac{1}{A_n(z)} \sim AR(n)$  in fact represents a stable

*AR* filter of order  $n$ , whose input error signal  $\varepsilon_n$  is white noise of constant spectral height equal to  $|T_n| / |T_{n-1}|$  and output is  $X_{n+1}$ . It can be shown that  $A_n(z)$  has all its zeros in  $|z| > 1$  provided  $|T_n| > 0$  thus establishing stability.

## Linear prediction Error

From (16-59), the mean square error using  $n$  samples is given by

$$\sigma_n^2 = \frac{|T_n|}{|T_{n-1}|} > 0. \quad (16-61)$$

Suppose one more sample from the past is available to evaluate  $X_{n+1}$  (i.e.,  $X_n, X_{n-1}, \dots, X_1, X_0$  are available). Proceeding as above the new coefficients and the mean square error  $\sigma_{n+1}^2$  can be determined. From (16-59)-(16-61),

$$\sigma_{n+1}^2 = \frac{|T_{n+1}|}{|T_n|}. \quad (16-62)$$

Using another matrix identity it is easy to show that

$$|T_{n+1}| = \frac{|T_n|^2}{|T_{n-1}|} (1 - |s_{n+1}|^2). \quad (16-63)$$

Since  $|T_k| > 0$ , we must have  $(1 - |s_{n+1}|^2) > 0$  or  $|s_{n+1}| < 1$  for every  $n$ . From (16-63), we have

$$\frac{|T_{n+1}|}{\underbrace{|T_n|}_{\sigma_{n+1}^2}} = \frac{|T_n|}{\underbrace{|T_{n-1}|}_{\sigma_n^2}} (1 - |s_{n+1}|^2)$$

or

$$\sigma_{n+1}^2 = \sigma_n^2 (1 - |s_{n+1}|^2) < \sigma_n^2, \quad (16-64)$$

since  $(1 - |s_{n+1}|^2) < 1$ . Thus the mean square error decreases as more and more samples are used from the past in the linear predictor.

In general from (16-64), the mean square errors for the one-step predictor form a monotonic nonincreasing sequence

$$\sigma_n^2 \geq \sigma_{n+1}^2 \geq \cdots \sigma_k^2 > \cdots \rightarrow \sigma_\infty^2 \quad (16-65)$$

whose limiting value  $\sigma_\infty^2 \geq 0$ .

Clearly,  $\sigma_\infty^2 \geq 0$  corresponds to the irreducible error in linear prediction using the entire past samples, and it is related to the power spectrum of the underlying process  $X(nT)$  through the relation

$$\sigma_\infty^2 = \exp \left[ \frac{1}{2\pi} \int_{-\pi}^{+\pi} \ln S_{XX}(\omega) d\omega \right] \geq 0. \quad (16-66)$$

where  $S_{XX}(\omega) \geq 0$  represents the power spectrum of  $X(nT)$ . For any finite power process, we have

$$\int_{-\pi}^{+\pi} S_{XX}(\omega) d\omega < \infty,$$

and since  $(S_{XX}(\omega) \geq 0)$ ,  $\ln S_{XX}(\omega) \leq S_{XX}(\omega)$ . Thus

$$\int_{-\pi}^{+\pi} \ln S_{XX}(\omega) d\omega \leq \int_{-\pi}^{+\pi} S_{XX}(\omega) d\omega < \infty. \quad (16-67)$$

Moreover, if the power spectrum is strictly positive at every Frequency, i.e.,

$$S_{xx}(\omega) > 0, \quad \text{in } -\pi < \omega < \pi, \quad (16-68)$$

then from (16-66)

$$\int_{-\pi}^{+\pi} \ln S_{xx}(\omega) d\omega > -\infty. \quad (16-69)$$

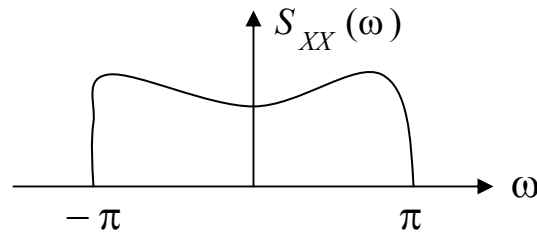
and hence

$$\sigma_{\infty}^2 = \exp \left[ \frac{1}{2\pi} \int_{-\pi}^{+\pi} \ln S_{xx}(\omega) d\omega \right] > e^{-\infty} = 0 \quad (16-70)$$

i.e., For processes that satisfy the strict positivity condition in (16-68) almost everywhere in the interval  $(-\pi, \pi)$ , the final minimum mean square error is strictly positive (see (16-70)).

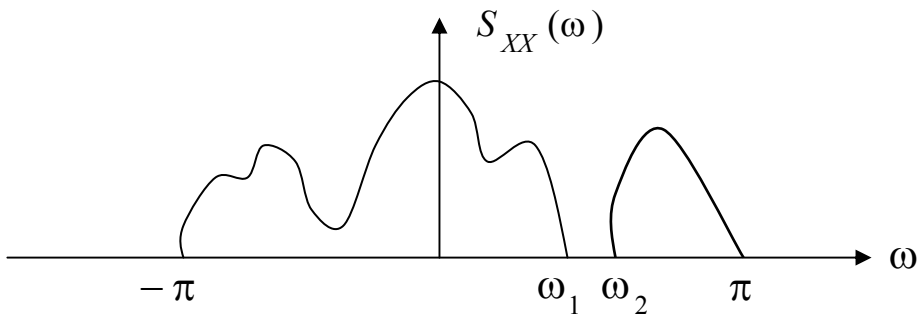
i.e., Such processes are not completely predictable even using their entire set of past samples, or they are inherently stochastic,

since the next output contains information that is not contained in the past samples. Such processes are known as *regular* stochastic processes, and their power spectrum is strictly positive.



Power Spectrum of a regular stochastic Process

Conversely, if a process has the following power spectrum,



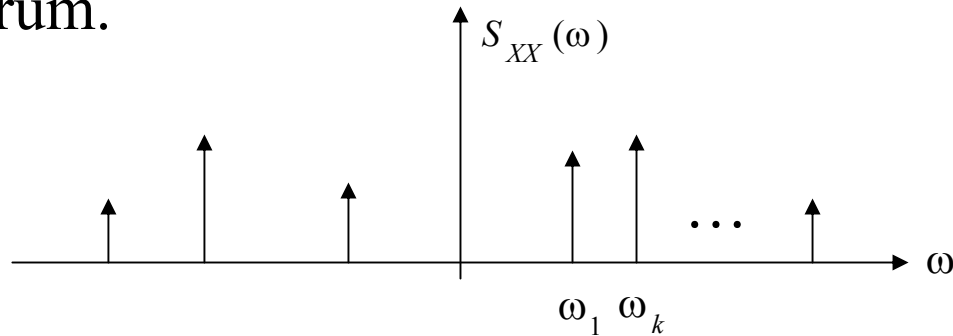
such that  $S_{XX}(\omega) = 0$  in  $\omega_1 < \omega < \omega_2$  then from (16-70),  $\sigma_{\infty}^2 = 0$ .



Such processes are completely predictable from their past data samples. In particular

$$X(nT) = \sum_k a_k \cos(\omega_k t + \phi_k) \quad (16-71)$$

is completely predictable from its past samples, since  $S_{XX}(\omega)$  consists of line spectrum.



$X(nT)$  in (16-71) is a shape deterministic stochastic process.