

VIII. Statistical Signal Processing

In this chapter we describe principles of statistical signal processing and tools that can be used to perform filtering, estimation, prediction and detection of random signals. We start with a review of discrete-time random processes and describe how signals and systems that occur in a variety of applications can be modelled by probabilistic tools. We then introduce signal modelling techniques which show how physical phenomena and systems can be represented by transfer functions, or equivalently, by filters or parameter vectors.

We then introduce the discrete-time Wiener filter which is a very powerful tool to estimate parameters, extract information and design signal processing systems. Linear prediction is then described along with the problem of performing prediction of signals and general parameters.

The chapter concludes with an introduction of adaptive signal processing and how a designer can employ adaptive algorithms to perform filtering, estimation, prediction and detection in the absence of statistical knowledge about signals and systems. We show that adaptive algorithms can learn the statistics about the environment and track slow variations.

A. Review of Random Processes

A random signal or a stochastic signal needs to be modelled using statistical information about it. Examples of such signals include speech, music and seismic signals. Discrete-time random signals are models of random processes, which are collections of random variables $X[n]$ characterized by a set of probability distribution functions. where n is the time index. In practice, each sample value $x[n]$ of a random signal is assumed to have resulted from a mechanism that is governed by a probability law.

Random variables:

A random variable $X[n]$ is described by the probability distribution function given by

$$P_{X[n]}(\alpha) = \Pr [X[n] < \alpha],$$

where α is a particular value of $X[n]$.

If $X[n]$ takes on a continuous range of values, it can be specified by the probability density function of $X[n]$:

$$p_{X[n]}(\alpha) = \frac{\partial P_{X[n]}(\alpha)}{\partial \alpha}$$

or equivalently

$$P_{X[n]}(\alpha) = \int_{-\infty}^{\alpha} p_{X[n]}(u) du$$

Properties of the probability density function:

$$p_{X[n]}(u) \geq 0$$

$$\int_{-\infty}^{\infty} p_{X[n]}(\alpha) d\alpha = 1$$

Properties of the probability distribution function:

$$0 \leq P_{X(n)}(\alpha) \leq 1$$

$$P_{X(n)}(\alpha_1) \leq P_{X(n)}(\alpha_2) \quad , \quad \text{for all } \alpha_2 > \alpha_1$$

$$P_{X(n)}(-\infty) = 0 \quad , \quad P_{X(n)}(+\infty) = 1$$

$$P_{X(n)}[\alpha_1 < X(n) \leq \alpha_2] = P_{X(n)}(\alpha_2) - P_{X(n)}(\alpha_1)$$

A random variable is characterized by a number of statistical properties. The r th moment of a random variable is described by

$$\mu_r = E[X(n)^r] = \int_{-\infty}^{\infty} \alpha^r P_{X(n)}(\alpha) d\alpha,$$

where r is any nonnegative integer and $E[\cdot]$ denotes the expectation operator. A random variable is completely characterized by all its moments.

The most commonly used statistical properties of a random variable are the mean or expected value m_x , the mean-square value $E[X(n)^2]$ and the variance σ_x^2 defined by:

$$m_x = E[X(n)] = \int_{-\infty}^{\infty} \alpha P_{X(n)}(\alpha) d\alpha$$

$$E[X(n)^2] = \int_{-\infty}^{\infty} \alpha^2 P_{X(n)}(\alpha) d\alpha$$

$$\sigma_x^2 = E[(X(n) - m_x)^2] = \int_{-\infty}^{\infty} (\alpha - m_x)^2 P_{X(n)}(\alpha) d\alpha.$$

There is also a relation between the quantities that is useful and is described by

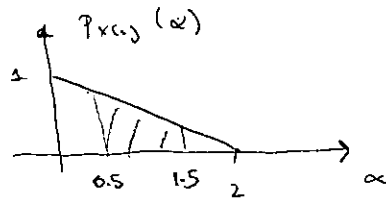
$$\sigma_x^2 = E[X(n)^2] - (m_x)^2.$$

where it turns out that $\sigma_x^2 = E[X(n)^2]$ for a random variable with zero mean.

Ex: Let a random variable be characterized by a probability density function

$$p_{X(n)}(\alpha) = \begin{cases} 1 - \frac{\alpha}{2}, & \text{for } 0 \leq \alpha \leq 2 \\ 0, & \text{otherwise} \end{cases}$$

as shown in the figure below



a) Compute its probability distribution function

$$\begin{aligned} P_{X(n)}(\alpha) &= \int_{-\infty}^{\alpha} p_{X(n)}(u) du \\ &= \int_0^{\alpha} \left(1 - \frac{u}{2}\right) du = \alpha - \frac{\alpha^2}{4}, \quad 0 < \alpha \leq 2 \end{aligned}$$

b) Compute the probability that $X(n)$ is in the range $0.5 < X(n) \leq 1.5$.

$$\begin{aligned} \Pr[0.5 < X(n) \leq 1.5] &= P_{X(n)}(1.5) - P_{X(n)}(0.5) \\ &= \frac{15}{16} - \frac{7}{16} = \frac{1}{2} \end{aligned}$$

c) Compute the mean, the mean-square value and the variance

$$\begin{aligned} m_x &= \int_0^2 \alpha \left(1 - \frac{\alpha}{2}\right) d\alpha = \int_0^2 \alpha - \frac{\alpha^2}{2} d\alpha \\ &= \left. \frac{\alpha^2}{2} - \frac{\alpha^3}{6} \right|_0^2 = \frac{4}{2} - \frac{8}{6} = \frac{4}{6} = \frac{2}{3} \end{aligned}$$

$$E[X(n)^2] = \int_0^2 \alpha^2 \left(1 - \frac{\alpha}{2}\right) d\alpha = \frac{2}{3}$$

$$\sigma_x^2 = E[X(n)^2] - (m_x)^2 = \frac{2}{3} - \left(\frac{2}{3}\right)^2 = \frac{2}{9}$$

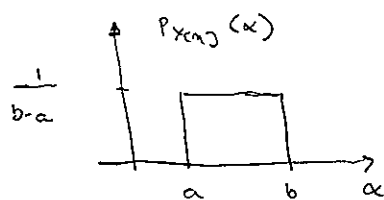
Amongst the most commonly used probability density functions in digital signal processing applications are the uniform density function and the Gaussian density function. The uniform density function is defined by

$$P_{X[n]}(\alpha) = \begin{cases} \frac{1}{b-a}, & \text{for } a \leq \alpha \leq b \\ 0, & \text{otherwise} \end{cases}$$

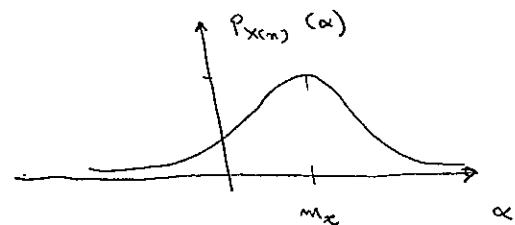
and the Gaussian density function is defined by

$$P_{X[n]}(\alpha) = \frac{1}{\sigma_x \sqrt{2\pi}} e^{-\frac{(\alpha - m_x)^2}{2\sigma_x^2}},$$

where the parameters m_x and σ_x are, respectively, the mean value and the standard deviation of $X[n]$ which lie in the range $-\infty < m_x < \infty$ and $\sigma_x > 0$. These density functions are illustrated in Fig. 1 below.



(a)



(b)

Fig. 1. (a) uniform and (b) Gaussian probability density functions.

Ex: Determine the mean and the variance of a uniformly distributed random variable.

$$\begin{aligned} m_x &= \int_a^b \alpha \cdot P_{X[n]}(\alpha) d\alpha = \int_a^b \alpha \cdot \frac{1}{b-a} d\alpha \\ &= \frac{1}{b-a} \left. \frac{\alpha^2}{2} \right|_a^b = \frac{1}{b-a} \frac{(b^2 - a^2)}{2} = \frac{b+a}{2} \end{aligned}$$

$$E[X^2[n]] = \frac{1}{b-a} \int_a^b \alpha^2 d\alpha = \frac{1}{b-a} \left. \frac{\alpha^3}{3} \right|_a^b = \frac{b^2 + a^2 + ab}{3}$$

$$\sigma_x^2 = E[X^2[n]] - (m_x)^2 = \frac{b^2 + a^2 + ab}{3} - \frac{b^2 + 2ab + a^2}{4} = \frac{(b-a)^2}{12}$$

When we deal with two random variables, say X and Y , it is also of interest to know their joint statistical properties as well as their individual statistical properties. The probability that $X(n)$ takes a value in a specified range from $-\infty$ to α and that $Y(n)$ takes a value in a specified range from $-\infty$ to β is given by their joint probability distribution function described by

$$P_{X(n)Y(n)}(\alpha, \beta) = \text{Pr} [X(n) \leq \alpha, Y(n) \leq \beta]$$

or equivalently, by their joint probability density function

$$f_{X(n)Y(n)}(\alpha, \beta) = \frac{\partial^2 P_{X(n)Y(n)}(\alpha, \beta)}{\partial \alpha \partial \beta}$$

The joint probability distribution function is thus given by

$$P_{X(n)Y(n)}(\alpha, \beta) = \int_{-\infty}^{\alpha} \int_{-\infty}^{\beta} f_{X(n)Y(n)}(u, v) du dv$$

Properties of the joint probability density function:

$$f_{X(n)Y(n)}(\alpha, \beta) \geq 0$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X(n)Y(n)}(\alpha, \beta) d\alpha d\beta = 1$$

Properties of the joint probability distribution function:

$$0 \leq P_{X(n)Y(n)}(\alpha, \beta) \leq 1$$

$$P_{X(n)Y(n)}(\alpha_1, \beta_1) \leq P_{X(n)Y(n)}(\alpha_2, \beta_2) \quad \text{for } \alpha_2 \geq \alpha_1 \text{ and } \beta_2 \geq \beta_1$$

$$P_{X(n)Y(n)}(-\infty, -\infty) = 0, \quad P_{X(n)Y(n)}(+\infty, +\infty) = 1$$

The joint statistical properties of two random variables $X(n)$ and $Y(n)$ are described by their cross-correlation and cross-covariance given by

$$r_{xy} = E[X(n)Y(n)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \alpha \beta f_{X(n)Y(n)}(\alpha, \beta) d\alpha d\beta$$

$$\begin{aligned} c_{xy} &= E[(X(n) - m_x)(Y(n) - m_y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\alpha - m_x) \cdot (\beta - m_y) f_{X(n)Y(n)}(\alpha, \beta) d\alpha d\beta \\ &= r_{xy} - m_x m_y \end{aligned}$$

where m_x and m_y are the mean of the random variables $X(n)$ and $Y(n)$.

Two random variables $X(n)$ and $Y(n)$ are said to be linearly independent or uncorrelated if

$$E[X(n)Y(n)] = E[X(n)]E[Y(n)]$$

and statistically independent if

$$P_{X(n), Y(n)}(\alpha, \beta) = P_{X(n)}(\alpha) \cdot P_{Y(n)}(\beta).$$

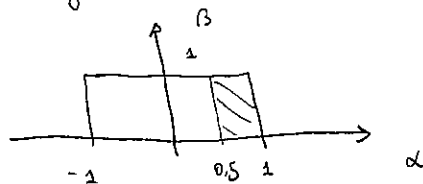
It can be shown that if two random variables are statistically independent then they are also uncorrelated. However, if these variables are uncorrelated they may not be statistically independent.

The statistical independence property makes it easier to compute the statistical properties of a random variable that is a function of several independent random variables

Ex: Consider the two random variables $X(n)$ and $Y(n)$ described by a uniformly distributed joint probability density function given by

$$P_{X(n), Y(n)}(\alpha, \beta) = \begin{cases} A, & -1 \leq \alpha \leq 1, 0 \leq \beta \leq 1, \\ 0, & \text{otherwise} \end{cases}$$

Determine the value of the constant A and then compute the probability that $X(n)$ and $Y(n)$ lie in the range $0.5 \leq \alpha \leq 1$, $0 \leq \beta \leq 1$ shown by



From the properties of joint probability density functions, we have

$$A \int_{-1}^1 \int_0^1 d\alpha d\beta = A \left[\int_{-1}^1 d\alpha \right] \left[\int_0^1 d\beta \right] = 2 \cdot A = 1$$

Hence, $A = 1/2$

$$P_F [0.5 \leq X(n) \leq 1, 0 \leq Y(n) \leq 1] = \frac{1}{2} \int_0^1 \int_{0.5}^1 d\alpha d\beta = \frac{1}{4}.$$

Statistical properties of random signals

In general, the mean, mean-square value and variance of a random discrete-time signal are functions of the time index n and can be considered as sequences.

So far we have considered real-valued random variables and random signals. A generalization to complex-valued random variables and signals is straight forward. The n th sample of a complex-valued random signal $x[n]$ is described by

$$x[n] = x_r[n] + j x_i[n],$$

where $x_r[n]$ and $x_i[n]$ are real-valued sequences called the real and imaginary parts of $x[n]$, respectively.

The mean value of a complex random signal is given by

$$\begin{aligned} m_x &= E[x[n]] = E[x_r[n] + j x_i[n]] \\ &= m_{x_r} + j m_{x_i} \end{aligned}$$

Likewise, the variance σ_x^2 of $x[n]$ is described by

$$\sigma_x^2 = E[|x[n] - m_x|^2] = E[|x[n]|^2] - (m_x)^2$$

The statistical relation of the samples of a random discrete-time signal at two different time indices m and n can be of great interest. One such relation is the autocorrelation, which for a complex random discrete-time signal $x[n]$ is given by

$$\Gamma_{xx}[m, n] = E[x[m] x^*[n]],$$

where $*$ denotes complex conjugation.

Another important relation is the autocovariance of $x[n]$ given by

$$\begin{aligned} C_{xx}[m, n] &= E[(x[m] - m_{x_m})(x[n] - m_{x_n})^*] \\ &= \Gamma_{xx}[m, n] - m_{x_m} m_{x_n}^* \end{aligned}$$

The correlation between two different random discrete-time signals $x[n]$ and $y[n]$ is described by the cross-correlation function

$$\begin{aligned} r_{xy}[m, n] &= E[x[n]y^*[m]] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \alpha \beta^* p_{x[n], y[m]}(\alpha, m, \beta, n) d\alpha d\beta \end{aligned}$$

and the cross-covariance function

$$\begin{aligned} c_{xy}[m, n] &= E[(x[n] - m_{x[n]})(y[m] - m_{y[m]})^*] \\ &= r_{xy}[m, n] - m_{x[n]} m_{y[m]}^* \end{aligned}$$

where $p_{x[n], y[m]}(\alpha, m, \beta, n)$ is the joint probability density function of $x[n]$ and $y[m]$. Both the cross-correlation and the cross-covariance functions can also be considered as two-dimensional sequences. The two random signals $x[n]$ and $y[n]$ are uncorrelated if $c_{xy}[m, n] = 0$ for all values of m and n .

Ex: Consider the random signal $x[n] = A \cos(\omega_0 n + \phi)$ and the prob. density functions of A and ϕ described by

$$p_A(\alpha) = \begin{cases} \frac{1}{4} & , 0 \leq \alpha \leq 4 \\ 0 & , \text{otherwise} \end{cases}$$

$$p_\phi(\phi) = \begin{cases} \frac{1}{2\pi} & , 0 \leq \phi \leq 2\pi \\ 0 & , \text{otherwise} \end{cases}$$

Since A and ϕ are statistically independent, their joint prob. density function is

$$p_{A\phi}(\alpha, \phi) = \begin{cases} \frac{1}{8\pi} & , 0 \leq \alpha \leq 4, 0 \leq \phi \leq 2\pi \\ 0 & , \text{otherwise} \end{cases}$$

The mean of $x[n]$ is

$$\begin{aligned} m_{x[n]} &= \frac{1}{8\pi} \int_0^4 \int_0^{2\pi} \alpha \cos(\omega_0 n + \phi) d\alpha d\phi \\ &= \frac{1}{8\pi} \left(\int_0^4 \alpha d\alpha \right) \left(\int_0^{2\pi} \cos(\omega_0 n + \phi) d\phi \right) = 0 \end{aligned}$$

The mean-square value is

$$\begin{aligned} E[x^2[n]] &= \frac{1}{8\pi} \int_0^4 \int_0^{2\pi} \alpha^2 \cos^2(\omega_0 n + \phi) d\alpha d\phi \\ &= \frac{1}{8\pi} \left(\int_0^4 \alpha^2 d\alpha \right) \left(\int_0^{2\pi} \cos^2(\omega_0 n + \phi) d\phi \right) = \frac{8}{3} \end{aligned}$$

The autocorrelation function is

$$\begin{aligned} r_{xx}[m, n] &= E[x[n]x[n]] \\ &= \frac{1}{8\pi} \int_0^4 \alpha^2 d\alpha \int_0^{2\pi} \cos(\omega_0 n + \phi) \cos(\omega_0 n + \phi) d\phi \\ &= \frac{8}{3} \cos(\omega_0(m-n)) \end{aligned}$$

Wide-sense stationary Random Signals:

A class of random signals often encountered in digital signal processing applications are the so-called wide-sense stationary (WSS) random processes for which some of the key statistical properties are all independent of time. Specifically, for a WSS random process $X[n]$, the mean $E[X[n]]$ has the same constant value m_x for all n , and the auto correlation and the autocovariance functions depend only on the difference of l and n , i.e.,

$$m_x = E[X[n]], \text{ for all } n$$

$$r_{xx}[l] = r_{xx}[n+l, n] = E[X[n+l]X^*[n]], \text{ for all } n \text{ and } l,$$

$$\begin{aligned} c_{xx}[l] &= c_{xx}[n+l, n] = E[(X[n+l]-m_x)(X[n]-m_x)^*] \\ &= r_{xx}[l] - |m_x|^2, \text{ for all } n \text{ and } l. \end{aligned}$$

The mean-square value of a WSS random process $X[n]$ is given by

$$E[|X[n]|^2] = r_{xx}[0],$$

and the variance is given by

$$\sigma_x^2 = c_{xx}[0] = r_{xx}[0] - |m_x|^2$$

The cross-correlation and cross-covariance functions between two WSS random processes $X[n]$ and $Y[n]$ are given by

$$r_{xy}[l] = E[X[n+l]Y^*[n]],$$

$$\begin{aligned} c_{xy}[l] &= E[(X[n+l]-m_x)(Y[n]-m_y)^*] \\ &= r_{xy}[l] - m_x m_y^* \end{aligned}$$

Concept of Power in a Random Signal :

In order to compute the power associated with a random signal $x[n]$ we employ the following definition:

$$P_x = E \left[\lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{m=-N}^N |x[m]|^2 \right]$$

In most practical cases, the expectation and the summation can be interchanged, resulting in a simpler expression given by

$$P_x = \lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{m=-N}^N E[|x[m]|^2]$$

In addition, if the random signal has a constant mean-square value, then the above expression reduces to

$$P_x = E[|x[m]|^2] = r_{xx}[0] = \sigma_x^2 + |m_x|^2.$$

Ergodic Signals

In many practical situations, the random signal of interest cannot be described in terms of a simple analytical expression and some of its statistical properties must be estimated from a finite set of samples. Such an approach can lead to meaningful results if the ergodicity condition is satisfied. In particular, a stationary random signal is defined to be an ergodic signal if all its statistical properties can be estimated from a single realization of sufficiently large finite length.

For an ergodic signal, time averages equal ensemble averages via the expectation operator in the limit as the length of the realization goes to infinity. For example, for a complex-valued ergodic signal we can compute the mean value, variance and auto covariance as :

$$m_x = \lim_{M \rightarrow \infty} \frac{1}{2M+1} \sum_{m=-M}^M x[m].$$

$$\sigma_x^2 = \lim_{M \rightarrow \infty} \frac{1}{2M+1} \sum_{n=-M}^M (x[n] - m_x) (x[n] - m_x)^*$$

$$r_{xx}[k] = \lim_{M \rightarrow \infty} \frac{1}{2M+1} \sum_{n=-M}^M (x[n] - m_x) (x[n+k] - m_x)^*$$

In practice, we replace the limit with a finite sum to provide an estimate of the desired statistical properties, as described by

$$\hat{m}_x = \frac{1}{M+1} \sum_{n=0}^M x[n],$$

$$\hat{\sigma}_x^2 = \frac{1}{M+1} \sum_{n=0}^M (x[n] - \hat{m}_x) (x[n] - \hat{m}_x)^*,$$

$$\hat{r}_{xx}[k] = \frac{1}{M+1} \sum_{n=0}^M (x[n] - \hat{m}_x) (x[n+k] - \hat{m}_x)^*.$$

Discrete-time processing of random signals:

In many situations, we need to study the effect of processing a random discrete-time signal by an LTI discrete-time system. Specifically, we need to determine the statistical properties of the output signal $y[n]$ generated by a stable LTI system with an impulse response $h[n]$ when its input $x[n]$ is a realization of a WSS random process $X[n]$.

Consider the linear convolution of $x[n]$ and $h[n]$ as described by

$$y[n] = \sum_{k=-\infty}^{\infty} h[k] \cdot x[n-k]$$

Since the function of a random variable is also a random variable, it follows that the output $y[n]$ is also a sample sequence of an output random process, $Y[n]$.

Since the input $x[n]$ is a sample sequence of $X[n]$, its mean m_x is a constant independent of the time index n . The mean $E[y[n]]$ of $y[n]$ is given by

$$\begin{aligned} m_y &= E[y[n]] = E \left[\sum_{k=-\infty}^{\infty} h[k] \cdot x[n-k] \right] \\ &= \sum_{k=-\infty}^{\infty} h[k] E[x[n-k]] = m_x \sum_{k=-\infty}^{\infty} h[k] = m_x H(e^{j0}) \end{aligned}$$

which is a constant.

The autocorrelation function of $y[n]$ is given by

$$\begin{aligned} r_{yy}[m+l, m] &= E[y[m+l] y^*[m]] \\ &= E\left[\left(\sum_{i=-\infty}^{\infty} h^*[i] x[m+l-i]\right) \left(\sum_{k=-\infty}^{\infty} h^*[k] x[m-k]\right)^*\right] \\ &= \sum_{i=-\infty}^{\infty} h^*[i] \sum_{k=-\infty}^{\infty} h[k] E[x[m+l-i] x^*[m-k]] \\ &= \sum_{i=-\infty}^{\infty} h^*[i] \sum_{k=-\infty}^{\infty} h[k] r_{xx}[m+l-i, m-k] \end{aligned}$$

which depends only on the difference $l+k-i$ of the time indices $m+l-i$ and $m-k$, i.e.,

$$r_{xx}[m+l-i, m-k] = r_{xx}[l+k-i].$$

Substituting the above into $r_{yy}[m+l, m]$, we arrive at

$$r_{yy}[m+l, m] = \sum_{i=-\infty}^{\infty} h^*[i] \sum_{k=-\infty}^{\infty} h[k] r_{xx}[l+k-i] = r_{yy}[l],$$

where the output autocorrelation sequence depends on l . Thus, the output $y[n]$ is a sample sequence of a WSS random process.

Substituting $m = i-k$ in the previous equation we arrive at

$$\begin{aligned} r_{yy}[l] &= \sum_{m=-\infty}^{\infty} r_{xx}[l-m] \sum_{k=-\infty}^{\infty} h^*[k] h[m+k] \\ &= \sum_{m=-\infty}^{\infty} r_{xx}[l-m] r_{hh}[m], \end{aligned}$$

where $r_{hh}[m]$ is the autocorrelation sequence of $h[n]$.

The cross-correlation function between the output and the input sequences of the LTI sequences of the LTI system for a real-valued input is given by

$$\begin{aligned} r_{yx}[m+l, m] &= E[y[m+l] x^*[m]] \\ &= E\left[\sum_{i=-\infty}^{\infty} h^*[i] x[m+l-i] x^*[m]\right] \\ &= \sum_{i=-\infty}^{\infty} h^*[i] E[x[m+l-i] x^*[m]] \\ &= \sum_{i=-\infty}^{\infty} h^*[i] r_{xx}[l-i] = r_{yx}[l], \end{aligned}$$

which indicates that the cross-correlation sequence depends on l .

B. Statistical Models

In statistical signal processing we employ models to represent or describe systems or physical phenomena, such systems are often excited by random signals and can be modelled by digital filters, which were described in the previous chapters. In general, we employ the following types of linear models:

- i) Autoregressive (AR) models
- ii) Moving average (MA) models
- iii) Auto-regressive moving average (ARMA) models

which can be depicted by the figure below.

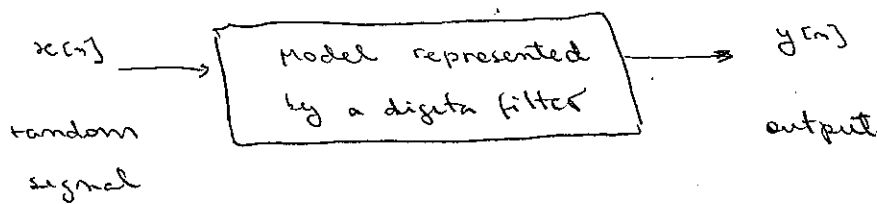


Fig. 2. Statistical model

Let us first consider an ARMA model which has an input $x[n]$ that corresponds to white noise and a rational transfer function given by

$$H(z) = \frac{B(z)}{A(z)} = \frac{\sum_{k=0}^q b[k] z^{-k}}{1 + \sum_{k=1}^p a[k] z^{-k}}$$

where $H(z)$ has p poles and q zeros, $b[k]$ and $a[k]$ are the parameters of the filter and it is assumed that the output of the filter $y[n]$ is WSS with $\sigma_{x[n]}^2 = \sigma_{y[n]}^2$. This filter corresponds to an IIR digital filter.

Since $x[n]$ and $y[n]$ are related by a difference equation given by

$$y[n] + \sum_{l=1}^p a[l] y[n-l] = \sum_{l=0}^q b[l] x[n-l]$$

it follows that the autocorrelation of $x[n]$ and the cross-correlation between $x[n]$ and $y[n]$ satisfy the same difference equation. Therefore, if we multiply both sides of the above equation by $y^*[n-k]$ and take the expected value, we have

$$r_{yy}[k] + \sum_{l=1}^p a[l] r_{yy}[k-l] = \sum_{l=0}^q b[l] \underbrace{E[x[n-l] y^*[n-k]]}_{r_{xy}[k-l]}$$

Therefore, the above eq. becomes

$$r_{yy}[k] + \sum_{l=1}^p a[l] r_{yy}[k-l] = \sum_{l=0}^q b[l] r_{xy}[k-l]$$

By writing $r_{xy}[k]$ in terms of $r_{yy}[k]$ and impulse response of the IIR filter $y[n] = h[n] * x[n] = \sum_{m=-\infty}^{\infty} h[m] \cdot x[n-m]$ the cross-correlation $r_{xy}[k-l]$ may be written as

$$\begin{aligned} E[x[n-l] \cdot y^*[n-k]] &= E\left[\sum_{m=-\infty}^{\infty} x[n-l] h^*[m] x[n-k-m]\right] \\ &= \sigma_x^2 h^*[l-k], \end{aligned}$$

where the last equality follows from $E[x[n-l] x^*[n-m]] = \sigma_x^2 \delta[n-l-m]$, i.e., $x[n]$ is white noise. Substituting the above eq into the previous one, we obtain

$$r_{yy}[k] + \sum_{l=1}^p a[l] r_{yy}[k-l] = \sigma_x^2 \sum_{l=0}^q b[l] \cdot h^*[l-k]$$

Assuming that $h[n]$ is causal, $\sigma_x^2 \sum_{l=0}^q b[l] \cdot h^*[l-k] = 0$

for $k > q$, we may write the above eq. as,

$$r_{yy}[k] + \sum_{l=1}^p a[l] \cdot r_{yy}[k-l] = \begin{cases} \sigma_x^2 \sum_{l=0}^q b[l] \cdot h^*[l-k], & 0 \leq k \leq q \\ 0, & k > q \end{cases}$$

which are the Yule-Walker equations

Writing the Yule-Walker equations for $k = 0, 1, \dots, p+q$ in matrix form, we have

$$\begin{bmatrix} r_{yy}[0] & r_{yy}[-1] & \dots & r_{yy}[-p] \\ r_{yy}[1] & r_{yy}[0] & \dots & r_{yy}[-p+1] \\ \vdots & \vdots & \ddots & \vdots \\ r_{yy}[q] & r_{yy}[q-1] & \dots & r_{yy}[q-p] \\ \hline r_{yy}[q+1] & r_{yy}[q] & \dots & r_{yy}[q-p+1] \\ \vdots & \vdots & \ddots & \vdots \\ r_{yy}[q+p] & r_{yy}[q+p-1] & \dots & r_{yy}[q] \end{bmatrix} \begin{bmatrix} 1 \\ a[1] \\ a[2] \\ \vdots \\ a[p] \end{bmatrix} = \Delta_n^2 \begin{bmatrix} \sum_{l=0}^q b[l] h^*[l] \\ \sum_{l=0}^{q-1} b[l] h^*[l-1] \\ \vdots \\ \sum_{l=0}^{q-p} b[l] h^*[l-p] \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

or

$$\underline{R} \underline{a} = \underline{\Gamma}$$

The above eq. defines a recursion for $r_{yy}[k]$ in terms of the filter coefficients $a[k]$ and $b[k]$. Therefore, the Yule-Walker equations may be used to compute $r_{yy}[k]$ from a finite set of $r_{yy}[k]$. For example, if $p > q$ and if $r_{yy}[0], \dots, r_{yy}[p-1]$ are known, then $r_{yy}[k]$ for $k \geq p$ may be computed as follows

$$r_{yy}[k] = - \sum_{l=1}^p a[l] \cdot r_{yy}[k-l]$$

The Yule-Walker equations may also be used to estimate the filter coefficients $a[k]$ and $b[k]$ from $r_{yy}[k]$. However, due to the product $h^*[l] b[k+l]$, the Yule-Walker equations are nonlinear in the filter coefficients and solving them for the filter coefficients is, in general, difficult.

The Yule-Walker equations are important in problems such as signal modelling and spectrum estimation. In what follows, we look at AR and MA models which are special cases of ARMA models.

In the AR model, $y[n]$ is generated by filtering $x[n]$ with an all-pole FIR filter of the form

$$H(z) = \frac{b[0]}{1 + \sum_{k=1}^p a[k] z^{-k}}$$

The Yule-Walker equations for an AR model can be found by setting $q=0$, which results in

$$r_{yy}[k] + \sum_{l=1}^p a[l] r_{yy}[k-l] = \sigma_x^2 |b[0]|^2 \delta[k], \quad k > 0$$

which in matrix form yields

$$\begin{bmatrix} r_{yy}[0] & r_{yy}[1] & \dots & r_{yy}[p] \\ r_{yy}[1] & r_{yy}[0] & \dots & r_{yy}[p-1] \\ \vdots & \vdots & \ddots & \vdots \\ r_{yy}[p] & r_{yy}[p-1] & \dots & r_{yy}[0] \end{bmatrix} \begin{bmatrix} 1 \\ a[1] \\ \vdots \\ a[p] \end{bmatrix} = \sigma_x^2 |b[0]|^2 \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Note that since the Yule-Walker equations are linear in $a[k]$, it is simple to find $a[k]$ from $r_{yy}[k]$

In the MA model, $y[n]$ is obtained by filtering $x[n]$ with an FIR filter that has a transfer function given by

$$H(z) = \sum_{k=0}^q b[k] z^{-k}$$

The Yule-Walker equations for an MA model can be found by setting $a[k]=0$ and noting that $h[n]=b[n]$, which yields

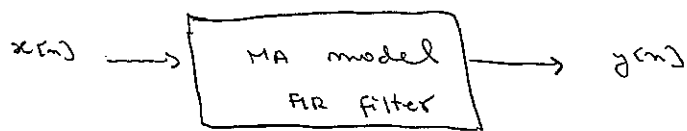
$$r_{yy}[k] = \sigma_x^2 \sum_{l=0}^{q-|k|} b[l+|k|] b^*[l]$$

Therefore, unlike the AR model the computation of an MA model to describe a process random signal is a difficult problem.

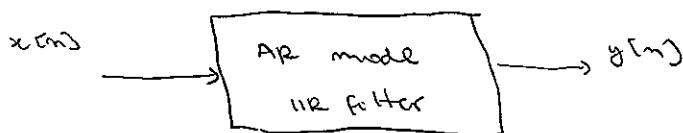
The use of AR models is widely advocated because of a result known as Wold's decomposition theorem. Such result states that a general linear model can be represented by an MA model which is equivalent to an AR model, i.e., whose coefficients obey the following relation

$$b[k] = a^*[k].$$

The basic difference between the MA and AR models is that $b[k]$ operates on the input $x[n]$ of the MA model, whereas $a[k]$ operates on the output of the AR model, as illustrated below.



(a)



(b)

Fig. 3 System models.

C. The Wiener Filter

In this section, we study an optimum filter called the Wiener filter that can be used to produce an optimum estimate of a signal from a noisy observation given by

$$x[n] = d[n] + w[n],$$

where $x[n]$ is a random signal that is observed at the input of a discrete-time LTI system, $d[n]$ is a desired signal and $w[n]$ represents noise. Both $d[n]$ and $w[n]$ are random signals taken from WSS random processes.

The problem we are interested in solving consists of designing an FIR filter $h[n]$ with order N that can result in the minimum mean-square error (MMSE) estimate of $d[n]$ by observing $x[n]$, as illustrated in the figure below.

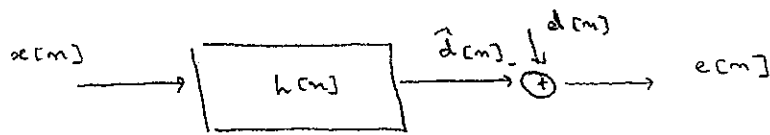


Fig. 4 The MMSE filtering problem.

The output $\hat{d}[n]$ of the filter $h[n]$ is described by

$$\hat{d}[n] = \sum_{k=0}^N h^*[k] x[n-k] = \underline{h}^H \underline{x}[n],$$

where the filter $h[n]$ and the input signal $x[n]$ can be represented using matrix notation; $\underline{h} = [h[0] \ h[1] \ \dots \ h[N]]^T$

$$\underline{x}[n] = [x[n] \ x[n-1] \ \dots \ x[n-N]]^T$$

The design problem corresponds to the following optimisation:

$$\begin{aligned} h_0 &= \arg \min C(h) = E[|e[m]|^2] \\ &= E\left[|d[m] - \underbrace{h^H x[m]}_{\hat{d}[m]}|^2\right] \end{aligned}$$

In order to compute the filter coefficients, we need to solve the minimisation of $C(h)$ by computing the derivative of $C(h)$ with respect to h^* or $h^*[k]$ and equating the result to zero as follows:

$$\begin{aligned} \frac{\partial C(h)}{\partial h^*[k]} &= \frac{\partial}{\partial h^*[k]} E[e[m] e^*[m]] \\ &= E\left[e[m] \frac{\partial}{\partial h^*[k]} e^*[m]\right] = 0, \quad k = 0, \dots, N \end{aligned}$$

where $e[m] = d[m] - \sum_{l=0}^N h^*[l] x[m-l]$ and $\frac{\partial}{\partial h^*[k]} e^*[m] = -x^*[m-k]$,

leading to the so-called orthogonality principle:

$$E[e^*[m] x[m-k]] = 0, \quad k = 0, 1, \dots, N$$

If we substitute the error $e[m]$ into the orthogonality principle, we can obtain the parameters for the Wiener filter as given by

$$\begin{aligned} E\left[\left(d[m] - \sum_{l=0}^N h^*[l] x[m-l]\right)^* x[m-k]\right] &= 0 \\ \underbrace{E[d^*[m] x[m-k]]}_{\Gamma_{dx}[k]} - \sum_{l=0}^N h^*[l] \underbrace{E[x^*[m-l] x[m-k]]}_{\Gamma_{xx}[k-l]} &= 0 \end{aligned}$$

$$\Gamma_{dx}[k] - \sum_{l=0}^N h^*[l] \Gamma_{xx}[k-l] = 0, \quad k = 0, 1, \dots, N,$$

where $\Gamma_{dx}[k]$ is the cross-correlation between $d[m]$ and $x[m-k]$ and $\Gamma_{xx}[k-l]$ is the autocorrelation of $x[m]$.

The Wiener filter $h[k]$ can also be employed and derived in matrix form. By substituting the previous expression for different values of k , we obtain

$$\sum_{l=0}^N h[l] r_{xx}[k-l] = r_{dx}[k], \quad k=0, 1, \dots, N$$

$$\begin{bmatrix} r_{xx}[0] & r_{xx}[1] & \dots & r_{xx}[N] \\ r_{xx}[1] & r_{xx}[0] & \dots & r_{xx}[N-1] \\ \vdots & \vdots & \ddots & \vdots \\ r_{xx}[N] & r_{xx}[N-1] & \dots & r_{xx}[0] \end{bmatrix} \begin{bmatrix} h[0] \\ h[1] \\ \vdots \\ h[N] \end{bmatrix} = \begin{bmatrix} r_{dx}[0] \\ r_{dx}[1] \\ \vdots \\ r_{dx}[N] \end{bmatrix},$$

$\underbrace{\hspace{15em}}_{\mathbf{R}_{xx}} \quad \underbrace{\hspace{5em}}_{\mathbf{h}_0} \quad \underbrace{\hspace{5em}}_{\mathbf{r}_{dx}}$

where \mathbf{R}_{xx} is the autocorrelation matrix, \mathbf{r}_{dx} is the cross-correlation vector and \mathbf{h}_0 is the Wiener filter.

The derivation of the Wiener filter can also be performed using matrix notation, which is more compact for obtaining the expression $\mathbf{R}_{xx} \mathbf{h}_0 = \mathbf{r}_{dx}$. Let us consider the optimisation problem given by

$$\mathbf{h}_0 = \arg \min_c C(\mathbf{h}) = E[|d[n] - \mathbf{h}^H \mathbf{z}[n]|^2]$$

This problem corresponds to computing the point of minimum of $C(\mathbf{h})$. In other words, it corresponds to calculating the derivative of $C(\mathbf{h})$ with respect to \mathbf{h} , equating the terms to zero and solving the resulting equations. To this end, it is convenient to define the partial derivatives of $C(\mathbf{h})$ with respect to \mathbf{h}^* using Wirtinger calculus:

$$\frac{\partial C}{\partial \mathbf{h}} = 0 \quad \text{and} \quad \frac{\partial C}{\partial \mathbf{h}^*} = \mathbf{0}$$

$$\begin{aligned}
\frac{\partial C(\underline{h})}{\partial \underline{h}^*} &= \frac{\partial}{\partial \underline{h}^*} E[|e[m]|^2] = 0 \\
&= \frac{\partial}{\partial \underline{h}^*} E[|d[m] - \underline{h}^H \underline{x}[m]|^2] = 0 \\
&= \frac{\partial}{\partial \underline{h}^*} E[(d[m] - \underline{h}^H \underline{x}[m])(d[m] - \underline{h}^H \underline{x}[m])^*] = 0 \\
&= E[-\underline{x}[m](d[m] - \underline{h}^H \underline{x}[m])^*] = 0
\end{aligned}$$

$$\underbrace{E[-\underline{x}[m] \cdot d[m]^*]}_{-\Gamma_{dx}} + \underbrace{E[\underline{x}[m] \underline{x}^H[m]]}_{R_{xx}} \underline{h}_0 = 0$$

$$R_{xx} \underline{h}_0 = \Gamma_{dx} \quad \rightarrow \quad \boxed{\underline{h}_0 = R_{xx}^{-1} \Gamma_{dx}}$$

The minimum mean-square error (MMSE) can be evaluated by substituting the Wiener filter \underline{h}_0 into the cost function as follows:

$$\begin{aligned}
\text{MMSE} = C(\underline{h}_0) &= E[|e[m]|^2] = E\left[e[m](d[m] - \sum_{l=0}^{N-1} \underline{h}_0[l] \underline{x}[m-l])^*\right] \\
&= E[e[m] d[m]^*] - \sum_{l=0}^{N-1} \underline{h}_0[l] E[e[m] \underline{x}^*[m-l]]
\end{aligned}$$

orthogonality principle.

Taking the expected values, we obtain

$$\begin{aligned}
\text{MMSE} &= E[d[m] d[m]^*] - \sum_{l=0}^{N-1} \underline{h}_0^*[l] \Gamma_{dx}[l] \\
&= \sigma_d^2 - \underline{h}_0^H \cdot \Gamma_{dx} \\
&= \sigma_d^2 - \Gamma_{dx}^H R_{xx}^{-1} \Gamma_{dx},
\end{aligned}$$

where the MMSE describes the accuracy of the Wiener filter to estimate $d[m]$.

Ex: Let $d[n]$ be an AR process with an autocorrelation sequence given by

$$r_{dx}[k] = \alpha^{|k|}$$

with $0 < \alpha < 1$ and suppose that $x[n]$ is observed in the presence of uncorrelated white noise that has a variance of σ_w^2 in the form:

$$x[n] = d[n] + w[n]$$

a) Design a Wiener filter to reduce the noise in $x[n]$

With $H(z) = \sum_{l=0}^N h[l] z^{-l}$ and $N=1$, we have

$$H(z) = h_0 + h_1 z^{-1}$$

The Wiener-Hopf equations are

$$\begin{bmatrix} r_{xx}[0] & r_{xx}[1] \\ r_{xx}[1] & r_{xx}[0] \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} = \begin{bmatrix} r_{dx}[0] \\ r_{dx}[1] \end{bmatrix}$$

Since $d[n]$ and $w[n]$ are assumed uncorrelated, we have

$$r_{dx}[k] = \alpha^{|k|}$$

$$\begin{aligned} r_{xx}[k] &= r_{dx}[k] + \underbrace{E[w[n]w^*[n+k]]}_{\sigma_w^2 = r_{ww}[k]} \\ &= r_{dx}[k] + \sigma_w^2 \end{aligned}$$

The Wiener-Hopf equations become

$$\begin{bmatrix} 1 + \sigma_w^2 & \alpha \\ \alpha & 1 + \sigma_w^2 \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \alpha \end{bmatrix}$$

Solving for h_0 and h_1 , we have

$$\begin{aligned} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} &= \frac{1}{(1 + \sigma_w^2) - \alpha^2} \begin{bmatrix} 1 + \sigma_w^2 & -\alpha \\ -\alpha & 1 + \sigma_w^2 \end{bmatrix} \begin{bmatrix} 1 \\ \alpha \end{bmatrix} \\ &= \frac{1}{(1 + \sigma_w^2) - \alpha^2} \begin{bmatrix} 1 + \sigma_w^2 & -\alpha^2 \\ \alpha & \alpha \sigma_w^2 \end{bmatrix} \end{aligned}$$

The Wiener filter is given by

$$H(z) = h_0 + h_1 z^{-1} = \frac{1}{(1+\alpha_w^2)^2 - \alpha^2} \left((1+\alpha_w^2 - \alpha^2) + \alpha \alpha_w^2 z^{-1} \right)$$

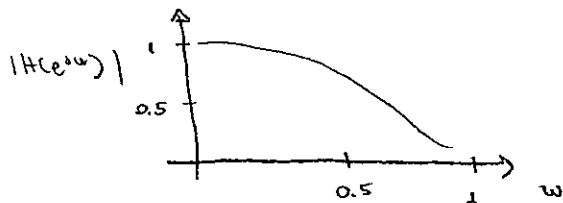
which has a zero at $z = \frac{-\alpha \alpha_w^2}{1+\alpha_w^2 - \alpha^2}$

b) Compute the power spectrum $H(e^{j\omega})$ for $\alpha = 0.8$ and $\alpha_w^2 = 1$

In this case, the Wiener filter becomes

$$H(e^{j\omega}) = 0.4048 + 0.2381 e^{-j\omega}$$

which is a lowpass filter with magnitude response



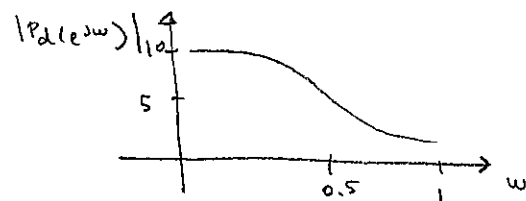
The fact that $H(e^{j\omega})$ is a lowpass filter is not a surprise.

The power spectrum of $d[n]$ is

$$P_d(e^{j\omega}) = \frac{1-\alpha^2}{(1+\alpha^2) - 2\alpha \cos \omega}$$

which for $\alpha = 0.8$ becomes

$$P_d(e^{j\omega}) = \frac{0.36}{1.64 - 1.6 \cos \omega}$$



Since $P_d(e^{j\omega})$ decreases with ω and the power spectrum of the noise is constant for all ω , then the signal-to-noise ratio decreases with the increase of ω .

c) Compute the MMSE for $\alpha = 0.8$ and $\alpha_w^2 = 1$

$$\begin{aligned} \text{MMSE} &= E[|e[n]|^2] = \sigma_d^2 - h^H \int d x \\ &= r_{dx}[0] - h_0^* r_{dx}[0] - h_1^* r_{dx}[1] \\ &= \alpha_w^2 \cdot \frac{1+\alpha_w^2 - \alpha^2}{(1+\alpha_w^2)^2 - \alpha^2} = 0.4048 \end{aligned}$$

d) Evaluate the increase in the signal-to-noise ratio (SNR) obtained by the Wiener filter for $\alpha = 0.8$ and $\sigma_w^2 = 1$.

Prior to filtering, we have

$$\text{SNR} = \frac{E[|d[m]|^2]}{E[|w[m]|^2]} = \frac{\sigma_d^2}{\sigma_w^2} = 1$$

$$\text{SNR}_{\text{dB}} = 10 \log_{10} 1 = 0 \text{ dB}$$

After filtering, we have

$$R_{xx}[k] = R_{dx}[k] + E_w[k]$$

The signal power is

$$E[|d'[m]|^2] = \mathbf{h}^H \cdot R_{xx} \mathbf{h} = [h_0^* \ h_1^*] \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} = 0.3748$$

The noise power is

$$E[|w'[m]|^2] = \mathbf{h}^H \cdot R_w \mathbf{h} = [h_0^* \ h_1^*] \begin{bmatrix} \sigma_w^2 & 0 \\ 0 & \sigma_w^2 \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} = 0.2206$$

The SNR at the output of the Wiener filter is

$$\text{SNR} = \frac{E[|d'[m]|^2]}{E[|w'[m]|^2]}$$

$$\text{SNR}_{\text{dB}} = 10 \log_{10} \frac{0.3748}{0.2206} = 2.302 \text{ dB}$$

Therefore, the Wiener filter increased the SNR by more than 2 dB.

D. Linear Prediction

Linear prediction is an important signal processing task which is encountered in many applications including communications, time-series analysis, control and biomedical engineering.

The basic problem is to design an FIR filter $h[n]$ to predict a desired signal $d[n+1]$ by observing $x[n] = d[n] + w[n]$. The input signal $x[n]$ contains a noisy version of the desired signal $d[n]$ in a previous time instant, as depicted in the figure below.

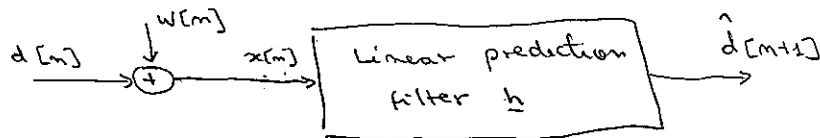


Fig. 5. Linear prediction problem.

Specifically, the linear prediction filter $h[n]$ performs linear prediction of $d[n]$ one step ahead. To this end, we consider the observed signal given by

$$x[n] = d[n] + w[n],$$

where both $d[n]$ and $w[n]$ are taken from WSS random processes.

The linear prediction filter design consists in designing an FIR filter $\underline{h} = [h[0] \ h[1] \ \dots \ h[N]]^T$ which will predict / estimate $d[n+1]$ based on the observation $\underline{x}[n] = [x[n] \ \dots \ x[n-N]]^T$. This is equivalent to the following linear combination:

$$\begin{aligned} \hat{d}[n+1] &= \sum_{k=0}^N h^*[k] x[n-k], \\ &= \underline{h}^H \cdot \underline{x}[n] \end{aligned}$$

The design of linear prediction filter \underline{h} requires the solution of the following optimisation problem

$$\begin{aligned} \underline{h}_0 &= \arg \min C_p(\underline{h}) = E \left[\underbrace{|d[m+1] - \hat{d}[m+1]|^2}_{e[m]} \right] \\ &= E \left[|d[m+1] - \sum_{k=0}^N h^*[k] x[m-k]|^2 \right] \end{aligned}$$

In order to compute the parameters of the linear prediction filter \underline{h} , we differentiate $C_p(\underline{h})$ with respect to \underline{h}^* , equate the terms to a null vector and solve the resulting equations:

$$\begin{aligned} \frac{\partial C_p(\underline{h})}{\partial \underline{h}^*} &= \frac{\partial}{\partial \underline{h}^*} E \left[(d[m+1] - \underline{h}^H x[m]) (d[m+1] - \underline{h}^H x[m])^* \right] = 0 \\ &= E \left[-x[m] (d[m+1] - \underline{h}^H x[m])^* \right] = 0 \\ &= -E \left[\underbrace{x[m] d^*[m+1]}_{\underline{\Gamma}_{dx}} \right] + E \left[\underbrace{x[m] x^H[m]}_{\underline{R}_{xx}} \right] \underline{h}_0 = 0 \end{aligned}$$

$$\underline{R}_{xx} \cdot \underline{h}_0 = \underline{\Gamma}_{dx} \Rightarrow \boxed{\underline{h}_0 = \underline{R}_{xx}^{-1} \underline{\Gamma}_{dx}},$$

where \underline{R}_{xx} is the autocorrelation matrix, $\underline{\Gamma}_{dx}$ is the cross-correlation vector and \underline{h}_0 is the optimal linear predictor. The linear prediction problem is similar to the Wiener filter problem and the main difference lies in the task performed by the filter.

The MMSE associated with the optimal linear predictor \underline{h}_0 is

$$\begin{aligned} \text{MMSE} &= C_p(\underline{h}_0) = E \left[|d[m+1] - \hat{d}[m+1]|^2 \right] \\ &= E \left[|d[m+1] - \sum_{k=0}^N h^*[k] x[m-k]|^2 \right] \\ &= E \left[e[m] d^*[m+1] - \sum_{k=0}^N h^*[k] E[e[m] x^*[m-k]] \right] \\ &= E \left[|d[m+1]|^2 \right] - \sum_{k=0}^N h^*[k] E \left[\underbrace{x[m-k] d^*[m+1]}_{\underline{\Gamma}_{dx}[k]} \right] \\ &= \sigma_d^2 - \underline{h}_0 \cdot \underline{\Gamma}_{dx}, \end{aligned}$$

where $\underline{\Gamma}_{dx}[k] = E[x[m-k] d^*[m+1]]$ is the cross-correlation function which is equivalent to $\underline{\Gamma}_{dx}[k+1]$.

Ex: Consider a linear prediction problem in which the observation is given by

$$x[m] = d[m] + w[m],$$

where $w[m]$ is a zero mean ($E[w[m]] = 0$) white noise with a variance of σ_w^2 , $d[m]$ and $w[m]$ are uncorrelated and $r_{dx}[k] = \alpha^{|k|}$.

a) Compute $r_{xx}[k]$ and $r_{dx}[k]$

$$r_{xx}[k] = E[x[m] x^*[m-k]] = r_{dx}[k] + r_{ww}[k] =$$

$$r_{dx}[k] = E[d[m+1] x^*[m-k]] = \alpha^{|k|}$$

$$r_{xx}[0] = 1 + \sigma_w^2, \quad r_{xx}[1] = \alpha$$

b) Compute the linear prediction filter

Since we know $r_{xx}[k]$ and $r_{dx}[k]$, we have

$$\underline{r}_{xx} \underline{h}_0 = \underline{r}_{dx} \Rightarrow \begin{bmatrix} 1 + \sigma_w^2 & \alpha \\ \alpha & 1 + \sigma_w^2 \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} = \begin{bmatrix} \alpha \\ \alpha^2 \end{bmatrix}$$

Solving for \underline{h}_0 , we obtain

$$\underline{h}_0 = \frac{\alpha}{(1 + \sigma_w^2) - \alpha^2} \begin{bmatrix} 1 + \sigma_w^2 - \alpha^2 \\ \alpha \sigma_w^2 \end{bmatrix}$$

E. Adaptive Algorithms

In the previous sections, we described the Wiener filter and the linear prediction filter. A key assumption was that the input signal $x[n]$ is WSS. In this case, we have shown how to design time-invariant FIR filters for filtering, estimation and prediction tasks, which result in an MMSE estimate of $d[n]$.

However, in many practical situations the signals that are processed are nonstationary and their statistics change with time. In order to deal with these signals, we may consider adaptive algorithms.

Let us consider the Wiener filtering problem and the input signal given by

$$x[n] = d[n] + w[n],$$

where $d[n]$ is nonstationary.

In this case, the Wiener filter is time-varying and produces the following estimate

$$\hat{d}[n] = \underline{h}^H[n] \cdot x[n]$$

In order to design the time-varying Wiener filter, we must solve the following optimisation problem

$$\min_{\underline{h}} C_{\pm}(\underline{h}) = E[|e[n]|^2] = E[|d[n] - \underline{h}^H[n] x[n]|^2]$$

The solution to this minimisation problem can be found by setting the derivative of $C_{\pm}(\underline{h})$ with respect to $\underline{h}^*[n]$ equal to zero.

This results in the following

$$E[x^*[n] x[n-k]] = 0, \quad k = 0, 1, \dots, N-1$$

Substituting the error $e[m] = d[m] - \underline{h}^H[m] \underline{x}[m]$ into the previous expression, we have

$$E \left[\left(d[m] - \sum_{\ell=0}^N h_{\ell}^*[m] x[m-\ell] \right)^* x[m-k] \right] = 0, \quad k=0, 1, \dots, N-1$$

which becomes

$$\sum_{\ell=0}^{N-1} h_{\ell}^*[m] E[x[m-\ell] x^*[m-k]] = E[d^*[m] x[m-k]], \quad k=0, 1, \dots, N-1$$

The solution to the above equation is the Wiener filter for nonstationary processes given by

$$\underline{h}[m] = \underline{R}_{xx}^{-1}[m] \underline{r}_{dx}[m],$$

where

$$\underline{R}_{xx}[m] = \begin{bmatrix} E[x[m]x^*[m]] & E[x[m-1]x^*[m]] & \dots & E[x[m-N]x^*[m]] \\ E[x[m]x^*[m-1]] & E[x[m-1]x^*[m-1]] & \dots & E[x[m-N]x^*[m-1]] \\ \vdots & \vdots & \ddots & \vdots \\ E[x[m]x^*[m-N]] & E[x[m-1]x^*[m-N]] & \dots & E[x[m-N]x^*[m-N]] \end{bmatrix}$$

$$\underline{r}_{dx}[m] = \begin{bmatrix} E[d^*[m]x[m]] \\ E[d^*[m]x[m-1]] \\ \vdots \\ E[d^*[m]x[m-N]] \end{bmatrix}$$

are the time-varying autocorrelation matrix and the cross-correlation vector, respectively.

The design of a time-varying filter is more difficult than a time-invariant one because the optimum set of coefficients must be computed every time instant n . The problem can be considerably simplified if we consider an approach that updates the coefficients in the form

$$\underline{h}[n+1] = \underline{h}[n] + \Delta \underline{h}[n],$$

where $\Delta \underline{h}[n]$ is a correction that is applied at time n to obtain the new set of coefficients at time $n+1$, as illustrated by the figure shown next.

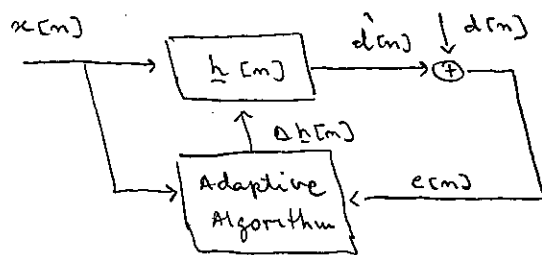


Fig. 6. Block diagram of an adaptive algorithm.

The fundamental approach that allows one to continuously update a set of parameters is at the heart of all adaptive signal processing algorithms. Each algorithm has its own way of obtaining the correction term $\Delta \hat{h}$. Moreover, the adaptive algorithm should have the following properties:

i) In a stationary environment, the algorithm should produce a sequence of corrections $\Delta \hat{h}$ in such a way that it converges to the Wiener solution, i.e.,

$$\lim_{n \rightarrow \infty} \hat{h}[n] = R_{xx}^{-1} P_{dx}$$

ii) The algorithm does not have to know the statistical quantities to compute $\Delta \hat{h}[n]$ as the algorithm should be able to compute these statistics with a built-in mechanism.

iii) For nonstationary signals, the filter should be able to adapt to the changing statistics and track the solution over time.

In what follows, we will describe the least-mean square (LMS) algorithm briefly along with a detailed derivation of it, conditions for its convergence and its requirements in terms of computational complexity.

The LMS algorithm was invented in 1960 by Widrow and Hopf and has become the most popular adaptive algorithm due to its simplicity, robustness and versatility. The LMS algorithm belongs to the class of stochastic gradient algorithms and employs instantaneous values of the gradient of $C_t(\underline{h}^H)$ with respect to \underline{h}^H in its adaptation.

In the derivation of the LMS algorithm, we consider the random input signal $x[n] = d[n] + w[n]$ and the filter $\underline{h}[n]$ with $N+1$ coefficients. The LMS solves the following optimisation problem

$$\min C_t(\underline{h}[n]) = E \left[\underbrace{|d[n] - \underline{h}^H[n] x[n]|^2}_{e[n]} \right]$$

By using a method known in optimisation as steepest descent and the gradient of $C_t(\underline{h}[n])$ with respect to $\underline{h}^H[n]$, we can search for the minimum of the cost function $C_t(\underline{h}[n])$ through the following recursion:

$$\begin{aligned} \underline{h}[n+1] &= \underline{h}[n] - \mu \nabla_{\underline{h}^H[n]} C_t(\underline{h}[n]) \\ &= \underline{h}[n] - \mu \left(\frac{\partial}{\partial \underline{h}^H[n]} E \left[(d[n] - \underline{h}^H[n] x[n]) \underbrace{(d[n] - \underline{h}^H[n] x[n])^*}_{e[n]} \right] \right) \\ &= \underline{h}[n] + \mu E \left[-x[n] d^*[n] + \underline{x}[n] x^H[n] \cdot \underline{h}[n] \right] \\ &= \underline{h}[n] + \mu \left(\underline{r}_d[n] - R_{xx}[n] \underline{h}[n] \right) \\ &= \underline{h}[n] + \mu E \left[e^*[n] \underline{x}[n] \right], \end{aligned}$$

where μ is the step size.

By employing instantaneous values of the gradient in the above development, we arrive at the LMS algorithm.

$$\begin{aligned} \underline{h}[n+1] &= \underline{h}[n] - \mu \nabla_{\underline{h}^H[n]} C_t(\underline{h}[n]) \\ &= \underline{h}[n] + \mu e^*[n] \cdot \underline{x}[n], \end{aligned}$$

where $e[n] = d[n] - \underline{h}^H[n] x[n]$ is the error signal.

The LMS algorithm has three basic steps, namely:

- filtering : $\hat{d}[n] = \mathbf{h}^H[n] \mathbf{x}[n]$
- Estimation error : $e[n] = d[n] - \hat{d}[n]$
- Adaptation : $\mathbf{h}[n+1] = \mathbf{h}[n] + \mu e^*[n] \mathbf{x}[n]$

These steps are illustrated in the block diagram below and in a table.

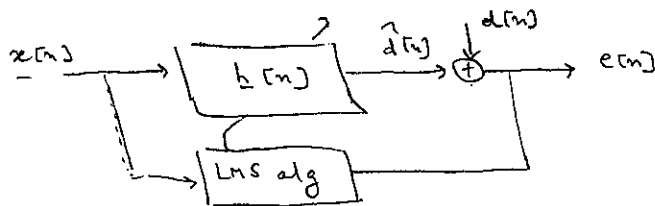


Fig. 7 Block diagram of the LMS algorithm

Table : The LMS algorithm

Initialisation:	
- N	filter length
- μ	step size
- $\mathbf{h}[0]$	= 0 - initial value of the filter
Computations:	
for $n = 1, 2, \dots$	
$\hat{d}[n]$	$= \mathbf{h}^H[n] \mathbf{x}[n]$
$e[n]$	$= d[n] - \hat{d}[n] = d[n] - \mathbf{h}^H[n] \mathbf{x}[n]$
$\mathbf{h}[n+1]$	$= \mathbf{h}[n] + \mu e^*[n] \mathbf{x}[n]$

In terms of computational complexity, the LMS algorithm has a cost that is linear in the number of filter coefficients, i.e., $O(N)$. It requires $2N+1$ complex multiplications and $2N$ complex additions, as shown below.

Table : Computational complexity of the LMS algorithm

Task	Multiplications	Additions
$\hat{d}[n] = \mathbf{h}^H[n] \mathbf{x}[n]$	N	$N-1$
$e[n] = d[n] - \hat{d}[n]$	-	1
$\mathbf{h}[n+1] = \mathbf{h}[n] + \mu e^*[n] \mathbf{x}[n]$	$N+1$	N
Total	$2N+1$	$2N$

The LMS algorithm works according to some general rules under which it converges and produces a given level of mean squared error (MSE). The LMS converges for the following conditions

$$0 < \mu < \frac{2}{\lambda_{\max}}$$

where λ_{\max} is the largest eigenvalue of $\mathbf{R}_{k \times k} = \sum_{k=0}^{N-1} \lambda_k \phi_k \phi_k^H$.

The MSE of the LMS algorithm can be evaluated analytically by the following expression.

$$\text{MSE}(m) = \text{MMSE} + \sum_{k=0}^{N-1} \frac{\lambda_k \cdot \mu}{2 - \mu \lambda_k} \text{MMSE}$$

which is illustrated by the following figure.

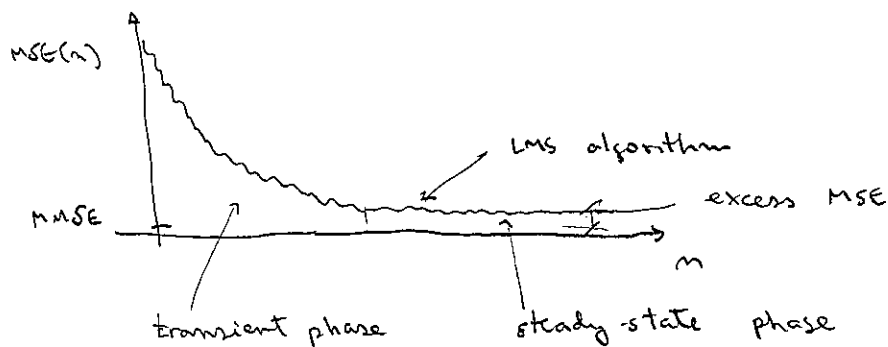


Fig. 8. Behaviour of the learning curve of the LMS algorithm.