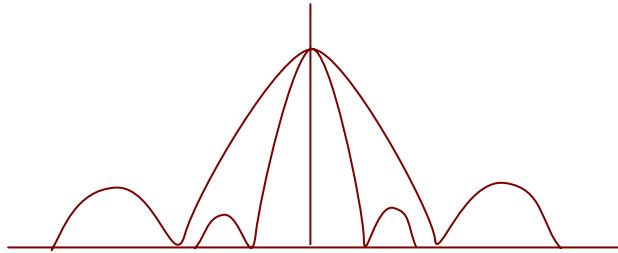


# **EC 622 Statistical Signal Processing**



**P. K. Bora**

**Department of Electronics & Communication Engineering  
INDIAN INSTITUTE OF TECHNOLOGY GUWAHATI**

## **EC 622 Statistical Signal Processing Syllabus**

1. Review of random variables: distribution and density functions, moments, independent, uncorrelated and orthogonal random variables; Vector-space representation of Random variables, Schwarz Inequality Orthogonality principle in estimation, Central Limit theorem, Random process, stationary process, autocorrelation and autocovariance functions, Spectral representation of random signals, Wiener Khinchin theorem, Properties of power spectral density, Gaussian Process and White noise process
2. Linear System with random input, Spectral factorization theorem and its importance, innovation process and whitening filter
3. Random signal modelling: MA(q), AR(p), ARMA(p,q) models
4. Parameter Estimation Theory: Principle of estimation and applications, Properties of estimates, unbiased and consistent estimators, MVUE, CR bound, Efficient estimators; Criteria of estimation: the methods of maximum likelihood and its properties; Bayesian estimation: Mean Square error and MMSE, Mean Absolute error, Hit and Miss cost function and MAP estimation
5. Estimation of signal in presence of White Gaussian Noise (WGN)  
  
Linear Minimum Mean-Square Error (LMMSE) Filtering: Wiener Hoff Equation  
FIR Wiener filter, Causal IIR Wiener filter, Noncausal IIR Wiener filter  
Linear Prediction of Signals, Forward and Backward Predictions, Levinson Durbin Algorithm, Lattice filter realization of prediction error filters
6. Adaptive Filtering: Principle and Application, Steepest Descent Algorithm  
Convergence characteristics; LMS algorithm, convergence, excess mean square error  
Leaky LMS algorithm; Application of Adaptive filters; RLS algorithm, derivation, Matrix inversion Lemma, Initialization, tracking of nonstationarity
7. Kalman filtering: Principle and application, Scalar Kalman filter, Vector Kalman filter
8. Spectral analysis: Estimated autocorrelation function, periodogram, Averaging the periodogram (Bartlett Method), Welch modification, Blackman and Tukey method of smoothing periodogram, Parametric method, AR(p) spectral estimation and detection of Harmonic signals, MUSIC algorithm.

## **Acknowledgement**

I take this opportunity to thank Prof. A. Mahanta who inspired me to take the course *Statistical Signal Processing*. I am also thankful to my other faculty colleagues of the ECE department for their constant support. I acknowledge the help of my students, particularly Mr. Diganta Gogoi and Mr. Gaurav Gupta for their help in preparation of the handouts. My appreciation goes to Mr. Sanjib Das who painstakingly edited the final manuscript and prepared the power-point presentations for the lectures. I acknowledge the help of Mr. L.N. Sharma and Mr. Nabajyoti Dutta for word-processing a part of the manuscript. Finally I acknowledge QIP, IIT Guwhati for the financial support for this work.

## **SECTION – I**

### **REVIEW OF RANDOM VARIABLES & RANDOM PROCESS**

## Table of Contents

<b>CHAPTER - 1: REVIEW OF RANDOM VARIABLES .....</b>	<b>9</b>
1.1 Introduction.....	9
1.2 Discrete and Continuous Random Variables .....	10
1.3 Probability Distribution Function .....	10
1.4 Probability Density Function .....	11
1.5 Joint random variable .....	12
1.6 Marginal density functions.....	12
1.7 Conditional density function.....	13
1.8 Baye's Rule for mixed random variables.....	14
1.9 Independent Random Variable .....	15
1.10 Moments of Random Variables .....	16
1.11 Uncorrelated random variables.....	17
1.12 Linear prediction of $Y$ from $X$ .....	17
1.13 Vector space Interpretation of Random Variables.....	18
1.14 Linear Independence .....	18
1.15 Statistical Independence.....	18
1.16 Inner Product .....	18
1.17 Schwarz Inequality .....	19
1.18 Orthogonal Random Variables.....	19
1.19 Orthogonality Principle.....	20
1.20 Chebysev Inequality.....	21
1.21 Markov Inequality .....	21
1.22 Convergence of a sequence of random variables .....	22
1.23 Almost sure (a.s.) convergence or convergence with probability 1 .....	22
1.24 Convergence in mean square sense .....	23
1.25 Convergence in probability .....	23
1.26 Convergence in distribution.....	24
1.27 Central Limit Theorem .....	24
1.28 Jointly Gaussian Random variables.....	25
<b>CHAPTER - 2 : REVIEW OF RANDOM PROCESS .....</b>	<b>26</b>
2.1 Introduction.....	26
2.2 How to describe a random process?.....	27
2.3 Stationary Random Process .....	28
2.4 Spectral Representation of a Random Process .....	30
2.5 Cross-correlation & Cross power Spectral Density.....	31
2.6 White noise process.....	32
2.7 White Noise Sequence.....	33
2.8 Linear Shift Invariant System with Random Inputs.....	33
2.9 Spectral factorization theorem .....	35
2.10 Wold's Decomposition .....	37
<b>CHAPTER - 3: RANDOM SIGNAL MODELLING .....</b>	<b>38</b>
3.1 Introduction.....	38
3.2 White Noise Sequence.....	38
3.3 Moving Average model $MA(q)$ model .....	38
3.4 Autoregressive Model .....	40
3.5 ARMA(p,q) – Autoregressive Moving Average Model .....	42
3.6 General $ARMA(p, q)$ Model building Steps .....	43
3.7 Other model: To model nonstationary random processes .....	43

<b>CHAPTER – 4: ESTIMATION THEORY .....</b>	<b>45</b>
4.1 Introduction .....	45
4.2 Properties of the Estimator .....	46
4.3 Unbiased estimator .....	46
4.4 Variance of the estimator .....	47
4.5 Mean square error of the estimator .....	48
4.6 Consistent Estimators .....	48
4.7 Sufficient Statistic .....	49
4.8 Cramer Rao theorem .....	50
4.9 Statement of the Cramer Rao theorem .....	51
4.10 Criteria for Estimation .....	54
4.11 Maximum Likelihood Estimator (MLE) .....	54
4.12 Bayesian Estimators .....	56
4.13 Bayesian Risk function or average cost .....	57
4.14 Relation between $\hat{\theta}_{\text{MAP}}$ and $\hat{\theta}_{\text{MLE}}$ .....	62
<b>CHAPTER – 5: WIENER FILTER.....</b>	<b>65</b>
5.1 Estimation of signal in presence of white Gaussian noise (WGN).....	65
5.2 Linear Minimum Mean Square Error Estimator .....	67
5.3 Wiener-Hopf Equations .....	68
5.4 FIR Wiener Filter .....	69
5.5 Minimum Mean Square Error - FIR Wiener Filter .....	70
5.6 IIR Wiener Filter (Causal).....	74
5.7 Mean Square Estimation Error – IIR Filter (Causal).....	76
5.8 IIR Wiener filter (Noncausal).....	78
5.9 Mean Square Estimation Error – IIR Filter (Noncausal).....	79
<b>CHAPTER – 6: LINEAR PREDICTION OF SIGNAL.....</b>	<b>82</b>
6.1 Introduction.....	82
6.2 Areas of application .....	82
6.3 Mean Square Prediction Error (MSPE) .....	83
6.4 Forward Prediction Problem .....	84
6.5 Backward Prediction Problem.....	84
6.6 Forward Prediction.....	84
6.7 Levinson Durbin Algorithm.....	86
6.8 Steps of the Levinson- Durbin algorithm.....	88
6.9 Lattice filter realization of Linear prediction error filters.....	89
6.10 Advantage of Lattice Structure .....	90
<b>CHAPTER – 7: ADAPTIVE FILTERS.....</b>	<b>92</b>
7.1 Introduction.....	92
7.2 Method of Steepest Descent.....	93
7.3 Convergence of the steepest descent method .....	95
7.4 Rate of Convergence .....	96
7.5 LMS algorithm (Least – Mean –Square) algorithm .....	96
7.6 Convergence of the LMS algorithm .....	99
7.7 Excess mean square error .....	100
7.8 Drawback of the LMS Algorithm.....	101
7.9 Leaky LMS Algorithm .....	103
7.10 Normalized LMS Algorithm .....	103
7.11 Discussion - LMS .....	104
7.12 Recursive Least Squares (RLS) Adaptive Filter .....	105

7.13 Recursive representation of $\hat{\mathbf{R}}_{YY}[n]$ .....	106
7.14 Matrix Inversion Lemma .....	106
7.15 RLS algorithm Steps .....	107
7.16 Discussion – RLS .....	108
7.16.1 Relation with Wiener filter .....	108
7.16.2. Dependence condition on the initial values .....	109
7.16.3. Convergence in stationary condition .....	109
7.16.4. Tracking non-staionarity .....	110
7.16.5. Computational Complexity .....	110
<b>CHAPTER – 8: KALMAN FILTER .....</b>	<b>111</b>
8.1 Introduction .....	111
8.2 Signal Model .....	111
8.3 Estimation of the filter-parameters .....	115
8.4 The Scalar Kalman filter algorithm .....	116
8.5 Vector Kalman Filter .....	117
<b>CHAPTER – 9 : SPECTRAL ESTIMATION TECHNIQUES FOR STATIONARY SIGNALS .....</b>	<b>119</b>
9.1 Introduction .....	119
9.2 Sample Autocorrelation Functions .....	120
9.3 Periodogram (Schuster, 1898) .....	121
9.4 Chi square distribution .....	124
9.5 Modified Periodograms .....	126
9.5.1 Averaged Periodogram: The Bartlett Method .....	126
9.5.2 Variance of the averaged periodogram .....	128
9.6 Smoothing the periodogram : The Blackman and Tukey Method .....	129
9.7 Parametric Method .....	130
9.8 AR spectral estimation .....	131
9.9 The Autocorrelation method .....	132
9.10 The Covariance method .....	132
9.11 Frequency Estimation of Harmonic signals .....	134
<b>10. Text and Reference .....</b>	<b>135</b>

## **SECTION – I**

### **REVIEW OF RANDOM VARIABLES & RANDOM PROCESS**



# CHAPTER - 1: REVIEW OF RANDOM VARIABLES

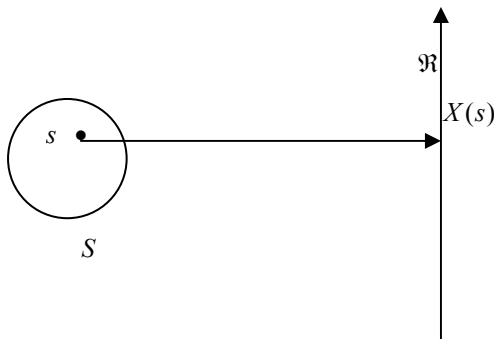
## 1.1 Introduction

- Mathematically a random variable is neither random nor a variable
- It is a mapping from sample space into the real-line ( “real-valued” random variable) or the complex plane ( “complex-valued ” random variable) .

Suppose we have a probability space  $\{S, \mathfrak{F}, P\}$  .

Let  $X : S \rightarrow \mathfrak{R}$  be a function mapping the sample space  $S$  into the real line such that For each  $s \in S$ , there exists a unique  $X(s) \in \mathfrak{R}$ . Then  $X$  is called a random variable.

Thus a random variable associates the points in the sample space with real numbers.



### Notations:

- Random variables are represented by upper-case letters.
- Values of a random variable are denoted by lower case letters
- $Y = y$  means that  $y$  is the value of a random variable  $X$ .

Figure Random Variable

**Example 1:** Consider the example of tossing a fair coin twice. The sample space is  $S = \{HH, HT, TH, TT\}$  and all four outcomes are equally likely. Then we can define a random variable  $X$  as follows

Sample Point	Value of the random Variable $X = x$	$P\{X = x\}$
HH	0	$\frac{1}{4}$
HT	1	$\frac{1}{4}$
TH	2	$\frac{1}{4}$
TT	3	$\frac{1}{4}$

**Example 2:** Consider the sample space associated with the single toss of a fair die. The sample space is given by  $S = \{1, 2, 3, 4, 5, 6\}$ . If we define the random variable  $X$  that associates a real number equal to the number in the face of the die, then  $X = \{1, 2, 3, 4, 5, 6\}$

## 1.2 Discrete and Continuous Random Variables

- A random variable  $X$  is called discrete if there exists a countable sequence of distinct real number  $x_i$  such that  $\sum_i P_m(x_i) = 1$ .  $P_m(x_i)$  is called the **probability mass function**. The random variable defined in Example 1 is a discrete random variable.
- A continuous random variable  $X$  can take any value from a continuous interval
- A random variable may also be mixed type. In this case the RV takes continuous values, but at each finite number of points there is a finite probability.

## 1.3 Probability Distribution Function

We can define an event  $\{X \leq x\} = \{s / X(s) \leq x, s \in S\}$

The probability

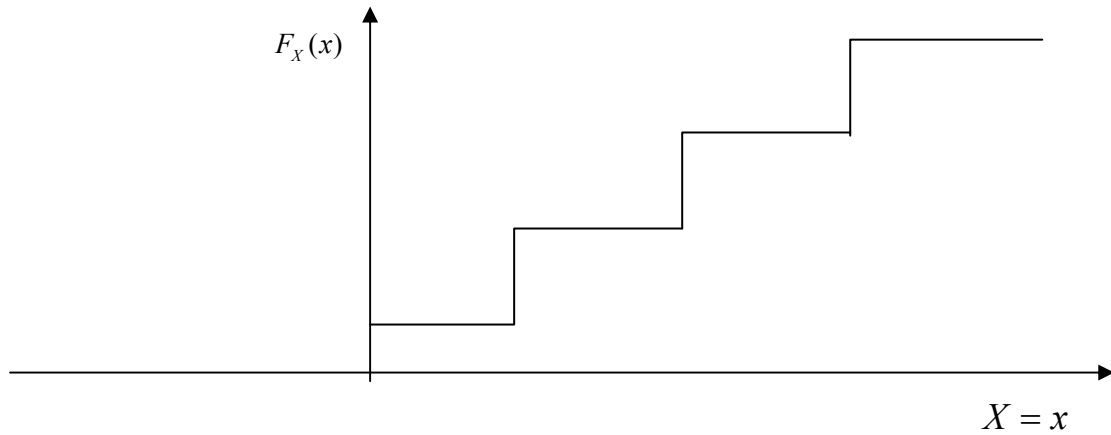
$F_X(x) = P\{X \leq x\}$  is called the probability distribution function.

Given  $F_X(x)$ , we can determine the probability of any event involving values of the random variable  $X$ .

- $F_X(x)$  is a non-decreasing function of  $X$ .
- $F_X(x)$  is right continuous  
 $\Rightarrow F_X(x)$  approaches to its value from right.
- $F_X(-\infty) = 0$
- $F_X(\infty) = 1$
- $P\{x_1 < X \leq x\} = F_X(x) - F_X(x_1)$

**Example 3:** Consider the random variable defined in Example 1. The distribution function  $F_X(x)$  is as given below:

Value of the random Variable $X = x$	$F_X(x)$
$x < 0$	0
$0 \leq x < 1$	$\frac{1}{4}$
$1 \leq x < 2$	$\frac{1}{2}$
$2 \leq x < 3$	$\frac{3}{4}$
$x \geq 3$	1



## 1.4 Probability Density Function

If  $F_X(x)$  is differentiable  $f_X(x) = \frac{d}{dx} F_X(x)$  is called the probability density function and has the following properties.

- $f_X(x)$  is a non-negative function
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$
- $P(x_1 < X \leq x_2) = \int_{-x_1}^{x_2} f_X(x) dx$

**Remark:** Using the Dirac delta function we can define the density function for a discrete random variables.

## 1.5 Joint random variable

$X$  and  $Y$  are two random variables defined on the same sample space  $S$ .

$P\{X \leq x, Y \leq y\}$  is called the joint distribution function and denoted by  $F_{X,Y}(x, y)$ .

Given  $F_{X,Y}(x, y)$ ,  $-\infty < x < \infty$ ,  $-\infty < y < \infty$ , we have a complete description of the random variables  $X$  and  $Y$ .

- $P\{0 < X \leq x, 0 < Y \leq y\} = F_{X,Y}(x, y) - F_{X,Y}(x, 0) - F_{X,Y}(0, y) + F_{X,Y}(0, 0)$
- $F_X(x) = F_{X,Y}(x, +\infty)$ .

To prove this

$$\begin{aligned} (X \leq x) &= (X \leq x) \cap (Y \leq +\infty) \\ \therefore F_X(x) &= P(X \leq x) = P(X \leq x, Y \leq \infty) = F_{X,Y}(x, +\infty) \\ F_X(x) &= P(X \leq x) = P(X \leq x, Y \leq \infty) = F_{X,Y}(x, +\infty) \end{aligned}$$

Similarly  $F_Y(y) = F_{X,Y}(\infty, y)$ .

- Given  $F_{X,Y}(x, y)$ ,  $-\infty < x < \infty$ ,  $-\infty < y < \infty$ , each of  $F_X(x)$  and  $F_Y(y)$  is called a marginal distribution function.

We can define joint probability density function  $f_{X,Y}(x, y)$  of the random variables  $X$  and  $Y$  by

$$f_{X,Y}(x, y) = \frac{\partial^2}{\partial x \partial y} F_{X,Y}(x, y), \text{ provided it exists}$$

- $f_{X,Y}(x, y)$  is always a positive quantity.
- $F_{X,Y}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{X,Y}(x, y) dx dy$

## 1.6 Marginal density functions

$$\begin{aligned} f_X(x) &= \frac{d}{dx} F_X(x) \\ &= \frac{d}{dx} F_{X,Y}(x, \infty) \\ &= \frac{d}{dx} \int_{-\infty}^x \left( \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \right) dx \\ &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \\ \text{and } f_Y(y) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx \end{aligned}$$

## 1.7 Conditional density function

$f_{Y/X}(y/X=x) = f_{Y/X}(y/x)$  is called conditional density of  $Y$  given  $X$ .

Let us define the conditional distribution function.

We cannot define the conditional distribution function for the continuous random variables  $X$  and  $Y$  by the relation

$$F_{Y/X}(y/x) = P(Y \leq y / X = x) \\ = \frac{P(Y \leq y, X = x)}{P(X = x)}$$

as both the numerator and the denominator are zero for the above expression.

The conditional distribution function is defined in the limiting sense as follows:

$$F_{Y/X}(y/x) = \lim_{\Delta x \rightarrow 0} P(Y \leq y / x < X \leq x + \Delta x) \\ = \lim_{\Delta x \rightarrow 0} \frac{P(Y \leq y, x < X \leq x + \Delta x)}{P(x < X \leq x + \Delta x)} \\ = \lim_{\Delta x \rightarrow 0} \frac{\int_{-\infty}^y f_{X,Y}(x, u) \Delta x du}{f_X(x) \Delta x} \\ = \frac{\int_{-\infty}^y f_{X,Y}(x, u) du}{f_X(x)}$$

The conditional density is defined in the limiting sense as follows

$$f_{Y/X}(y/X=x) = \lim_{\Delta y \rightarrow 0} (F_{Y/X}(y + \Delta y / X = x) - F_{Y/X}(y / X = x)) / \Delta y \\ = \lim_{\Delta y \rightarrow 0, \Delta x \rightarrow 0} (F_{Y/X}(y + \Delta y / x < X \leq x + \Delta x) - F_{Y/X}(y / x < X \leq x + \Delta x)) / \Delta y \quad (1)$$

Because  $P(X = x) = \lim_{\Delta x \rightarrow 0} P(x < X \leq x + \Delta x)$

The right hand side in equation (1) is

$$\lim_{\Delta y \rightarrow 0, \Delta x \rightarrow 0} (F_{Y/X}(y + \Delta y / x < X < x + \Delta x) - F_{Y/X}(y / x < X < x + \Delta x)) / \Delta y \\ = \lim_{\Delta y \rightarrow 0, \Delta x \rightarrow 0} (P(y < Y \leq y + \Delta y / x < X \leq x + \Delta x)) / \Delta y \\ = \lim_{\Delta y \rightarrow 0, \Delta x \rightarrow 0} (P(y < Y \leq y + \Delta y, x < X \leq x + \Delta x)) / P(x < X \leq x + \Delta x) \Delta y \\ = \lim_{\Delta y \rightarrow 0, \Delta x \rightarrow 0} f_{X,Y}(x, y) \Delta x \Delta y / f_X(x) \Delta x \Delta y \\ = f_{X,Y}(x, y) / f_X(x)$$

$$\boxed{\therefore f_{Y/X}(x/y) = f_{X,Y}(x, y) / f_X(x)} \quad (2)$$

Similarly, we have

$$\boxed{\therefore f_{X/Y}(x/y) = f_{X,Y}(x, y) / f_Y(y)} \quad (3)$$

From (2) and (3) we get Baye's rule

$$\begin{aligned}
 \therefore f_{X/Y}(x/y) &= \frac{f_{X,Y}(x,y)}{f_Y(y)} \\
 &= \frac{f_X(x)f_{Y/X}(y/x)}{f_Y(y)} \\
 &= \frac{f_{X,Y}(x,y)}{\int_{-\infty}^{\infty} f_{X,Y}(x,y)dx} \\
 &= \frac{f_{Y/X}(y/x)f_X(x)}{\int_{-\infty}^{\infty} f_X(u)f_{Y/X}(y/x)du}
 \end{aligned} \tag{4}$$

**Given the joint density function we can find out the conditional density function.**

#### **Example 4:**

For random variables  $X$  and  $Y$ , the joint probability density function is given by

$$\begin{aligned}
 f_{X,Y}(x,y) &= \frac{1+xy}{4} \quad |x| \leq 1, \quad |y| \leq 1 \\
 &= 0 \quad \text{otherwise}
 \end{aligned}$$

Find the marginal density  $f_X(x)$ ,  $f_Y(y)$  and  $f_{Y/X}(y/x)$ . Are  $X$  and  $Y$  independent?

$$f_X(x) = \int_{-1}^1 \frac{1+xy}{4} dy = \frac{1}{2}$$

Similarly

$$f_Y(y) = \frac{1}{2} \quad -1 \leq y \leq 1$$

and

$$\begin{aligned}
 f_{Y/X}(y/x) &= \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{1+xy}{4}, \quad |x| \leq 1, \quad |y| \leq 1 \\
 &= 0 \quad \text{otherwise}
 \end{aligned}$$

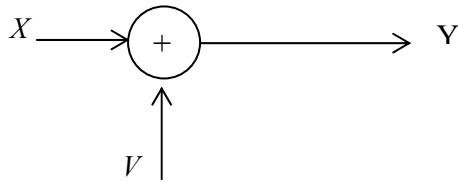
$\therefore X$  and  $Y$  are not independent

### **1.8 Baye's Rule for mixed random variables**

Let  $X$  be a discrete random variable with probability mass function  $P_X(x)$  and  $Y$  be a continuous random variable. In practical problem we may have to estimate  $X$  from observed  $Y$ . Then

$$\begin{aligned}
P_{X/Y}(x/y) &= \lim_{\Delta y \rightarrow 0} \frac{P_{X/Y}(x/y < Y \leq y + \Delta y)}{\Delta y} \\
&= \lim_{\Delta y \rightarrow 0} \frac{P_{X,Y}(x, y < Y \leq y + \Delta y)}{P_Y(y < Y \leq y + \Delta y)} \\
&= \lim_{\Delta y \rightarrow 0} \frac{P_X(x) f_{Y/X}(y/x) \Delta y}{f_Y(y) \Delta y} \\
&= \frac{P_X(x) f_{Y/X}(y/x)}{f_Y(y)} \\
&= \frac{P_X(x) f_{Y/X}(y/x)}{\sum_x P_X(x) f_{Y/X}(y/x)}
\end{aligned}$$

**Example 5:**



$X$  is a binary random variable with

$$X = \begin{cases} 1 & \text{with probability } \frac{1}{2} \\ -1 & \text{with probability } \frac{1}{2} \end{cases}$$

$V$  is the Gaussian noise with mean 0 and variance  $\sigma^2$ .

Then

$$\begin{aligned}
P_{X/Y}(x=1/y) &= \frac{P_X(x) f_{Y/X}(y/x)}{\sum_x P_X(x) f_{Y/X}(y/x)} \\
&= \frac{e^{-(y-1)^2/2\sigma^2}}{e^{-(y-1)^2/2\sigma^2} + e^{-(y+1)^2/2\sigma^2}}
\end{aligned}$$

## 1.9 Independent Random Variable

Let  $X$  and  $Y$  be two random variables characterised by the joint density function

$$F_{X,Y}(x,y) = P\{X \leq x, Y \leq y\}$$

$$\text{and } f_{X,Y}(x,y) = \frac{\partial^2}{\partial x \partial y} F_{X,Y}(x,y)$$

Then  $X$  and  $Y$  are independent if  $f_{X/Y}(x/y) = f_X(x) \quad \forall x \in \mathfrak{R}$

and equivalently

$f_{X,Y}(x,y) = f_X(x) f_Y(y)$  , where  $f_X(x)$  and  $f_Y(y)$  are called the marginal density functions.

## 1.10 Moments of Random Variables

- Expectation provides a description of the random variable in terms of a few parameters instead of specifying the entire distribution function or the density function
- It is far easier to estimate the expectation of a R.V. from data than to estimate its distribution

### First Moment or mean

The mean  $\mu_X$  of a random variable  $X$  is defined by

$$\begin{aligned}\mu_X &= EX = \sum x_i P(x_i) \text{ for a discrete random variable } X \\ &= \int_{-\infty}^{\infty} x f_X(x) dx \text{ for a continuous random variable } X\end{aligned}$$

For any piecewise continuous function  $y = g(x)$ , the expectation of the R.V.

$$Y = g(X) \text{ is given by } EY = Eg(X) = \int_{-\infty}^{\infty} g(x) f_X(x) dx$$

### Second moment

$$EX^2 = \int_{-\infty}^{\infty} x^2 f_X(x) dx$$

### Variance

$$\sigma_X^2 = \int_{-\infty}^{\infty} (x - \mu_X)^2 f_X(x) dx$$

- Variance is a central moment and measure of dispersion of the random variable about the mean.
- $\sigma_x$  is called the **standard deviation**.
- 

For two random variables  $X$  and  $Y$  the joint expectation is defined as

$$E(XY) = \mu_{X,Y} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{X,Y}(x,y) dx dy$$

The correlation between random variables  $X$  and  $Y$ , measured by the covariance, is given by

$$\begin{aligned}\text{Cov}(X,Y) &= \sigma_{XY} = E(X - \mu_X)(Y - \mu_Y) \\ &= E(XY - X\mu_Y - Y\mu_X + \mu_X\mu_Y) \\ &= E(XY) - \mu_X\mu_Y\end{aligned}$$

The ratio 
$$\rho = \frac{E(X - \mu_X)(Y - \mu_Y)}{\sqrt{E(X - \mu_X)^2 E(Y - \mu_Y)^2}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

is called the **correlation coefficient**. The correlation coefficient measures how much two random variables are similar.



## 1.11 Uncorrelated random variables

Random variables  $X$  and  $Y$  are uncorrelated if covariance

$$\text{Cov}(X, Y) = 0$$

Two random variables may be dependent, but still they may be uncorrelated. If there exists correlation between two random variables, one may be represented as a linear regression of the others. We will discuss this point in the next section.

## 1.12 Linear prediction of $Y$ from $X$

$$\hat{Y} = aX + b \quad \text{Regression}$$

$$\text{Prediction error } Y - \hat{Y}$$

Mean square prediction error

$$E(Y - \hat{Y})^2 = E(Y - aX - b)^2$$

For minimising the error will give optimal values of  $a$  and  $b$ . Corresponding to the optimal solutions for  $a$  and  $b$ , we have

$$\frac{\partial}{\partial a} E(Y - aX - b)^2 = 0$$

$$\frac{\partial}{\partial b} E(Y - aX - b)^2 = 0$$

$$\text{Solving for } a \text{ and } b, \quad \hat{Y} - \mu_Y = \frac{1}{\sigma_X^2} \sigma_{X,Y} (x - \mu_X)$$

$$\text{so that } \hat{Y} - \mu_Y = \rho_{X,Y} \frac{\sigma_Y}{\sigma_X} (x - \mu_X), \text{ where } \boxed{\rho_{X,Y} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}} \text{ is the correlation coefficient.}$$

If  $\rho_{X,Y} = 0$  then  $X$  and  $Y$  are uncorrelated.

$$\Rightarrow \hat{Y} - \mu_Y = 0$$

$$\Rightarrow \hat{Y} = \mu_Y \text{ is the best prediction.}$$

Note that independence  $\Rightarrow$  Uncorrelatedness. But uncorrelated generally does not imply independence (except for jointly Gaussian random variables).

### Example 6:

$$Y = X^2 \text{ and } f_X(x) \text{ is uniformly distributed between } (1, -1).$$

$X$  and  $Y$  are dependent, but they are uncorrelated.

$$\text{Cov}(X, Y) = \sigma_X = E(X - \mu_X)(Y - \mu_Y)$$

$$\begin{aligned} \text{Because} \quad &= EXY = EX^3 = 0 \\ &= EXEY \quad (\because EX = 0) \end{aligned}$$

In fact for any zero-mean symmetric distribution of  $X$ ,  $X$  and  $X^2$  are uncorrelated.

### 1.13 Vector space Interpretation of Random Variables

The set of all random variables defined on a **sample space** form a vector space with respect to addition and scalar multiplication. This is very easy to verify.

### 1.14 Linear Independence

Consider the sequence of random variables  $X_1, X_2, \dots, X_N$ .

If  $c_1X_1 + c_2X_2 + \dots + c_NX_N = 0$  implies that

$$c_1 = c_2 = \dots = c_N = 0, \text{ then } X_1, X_2, \dots, X_N \text{ are linearly independent.}$$

### 1.15 Statistical Independence

$X_1, X_2, \dots, X_N$  are statistically independent if

$$f_{X_1, X_2, \dots, X_N}(x_1, x_2, \dots, x_N) = f_{X_1}(x_1)f_{X_2}(x_2) \dots f_{X_N}(x_N)$$

Statistical independence in the case of zero mean random variables also implies linear independence

### 1.16 Inner Product

If  $x$  and  $y$  are real vectors in a vector space  $V$  defined over the field  $\mathbb{R}$ , the inner product  $\langle x, y \rangle$  is a scalar such that

$$\forall x, y, z \in V \text{ and } a \in \mathbb{R}$$

1.  $\langle x, y \rangle = \langle y, x \rangle$
2.  $\langle x, x \rangle = \|x\|^2 \geq 0$
3.  $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$
4.  $\langle ax, y \rangle = a \langle x, y \rangle$

In the case of RVs, inner product between  $X$  and  $Y$  is defined as

$$\langle X, Y \rangle = EXY = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{X,Y}(x, y) dy dx.$$

Magnitude / Norm of a vector

$$\|x\|^2 = \langle x, x \rangle$$

So, for R.V.

$$\|X\|^2 = EX^2 = \int_{-\infty}^{\infty} x^2 f_X(x) dx$$

- The set of RVs along with the inner product defined through the joint expectation operation and the corresponding norm defines a Hilbert Space.

## 1.17 Schwarz Inequality

For any two vectors  $x$  and  $y$  belonging to a Hilbert space  $V$

$$|\langle x, y \rangle| \leq \|x\| \|y\|$$

For RV  $X$  and  $Y$

$$E^2(XY) \leq EX^2 EY^2$$

### Proof:

Consider the random variable  $Z = aX + Y$

$$\begin{aligned} E(aX + Y)^2 &\geq 0 \\ \Rightarrow a^2 EX^2 + EY^2 + 2aEXY &\geq 0 \end{aligned}$$

Non-negativity of the left-hand side  $\Rightarrow$  its minimum also must be nonnegative.

For the minimum value,

$$\frac{dEZ^2}{da} = 0 \Rightarrow a = -\frac{EXY}{EX^2}$$

so the corresponding minimum is  $\frac{E^2 XY}{EX^2} + EY^2 - 2\frac{E^2 XY}{EX^2}$

Minimum is nonnegative  $\Rightarrow$

$$EY^2 - \frac{E^2 XY}{EX^2} \geq 0$$

$$\Rightarrow \boxed{E^2 XY < EX^2 EY^2}$$

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E(X - \mu_X)(Y - \mu_Y)}{\sqrt{E(X - \mu_X)^2 E(Y - \mu_Y)^2}}$$

From Schwarz inequality

$$\|\rho(X, Y)\| \leq 1$$

## 1.18 Orthogonal Random Variables

Recall the definition of orthogonality. Two vectors  $x$  and  $y$  are called orthogonal if

$$\langle x, y \rangle = 0$$

Similarly two random variables  $X$  and  $Y$  are called orthogonal if  $EXY = 0$

If each of  $X$  and  $Y$  is zero-mean

$$\text{Cov}(X, Y) = EXY$$

Therefore, if  $EXY = 0$  then  $\text{Cov}(X, Y) = 0$  for this case.

For zero-mean random variables,

$$\text{Orthogonality} \Leftrightarrow \text{uncorrelatedness}$$

## 1.19 Orthogonality Principle

$X$  is a random variable which is not observable.  $Y$  is another observable random variable which is statistically dependent on  $X$ . Given a value of  $Y$  what is the best guess for  $X$ ? (Estimation problem).

Let the best estimate be  $\hat{X}(Y)$ . Then  $E(X - \hat{X}(Y))^2$  is a minimum with respect to  $\hat{X}(Y)$ .

And the corresponding estimation principle is called *minimum mean square error principle*. For finding the minimum, we have

$$\begin{aligned}\frac{\partial}{\partial \hat{X}} E(X - \hat{X}(Y))^2 &= 0 \\ \Rightarrow \frac{\partial}{\partial \hat{X}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \hat{X}(y))^2 f_{X,Y}(x, y) dy dx &= 0 \\ \Rightarrow \frac{\partial}{\partial \hat{X}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \hat{X}(y))^2 f_Y(y) f_{X/Y}(x) dy dx &= 0 \\ \Rightarrow \frac{\partial}{\partial \hat{X}} \int_{-\infty}^{\infty} f_Y(y) \left( \int_{-\infty}^{\infty} (x - \hat{X}(y))^2 f_{X/Y}(x) dx \right) dy &= 0\end{aligned}$$

Since  $f_Y(y)$  in the above equation is always positive, therefore the minimization is equivalent to

$$\begin{aligned}\frac{\partial}{\partial \hat{X}} \int_{-\infty}^{\infty} (x - \hat{X}(y))^2 f_{X/Y}(x) dx &= 0 \\ \text{Or } 2 \int_{-\infty}^{\infty} (x - \hat{X}(y)) f_{X/Y}(x) dx &= 0 \\ \Rightarrow \int_{-\infty}^{\infty} \hat{X}(y) f_{X/Y}(x) dx &= \int_{-\infty}^{\infty} x f_{X/Y}(x) dx \\ \Rightarrow \hat{X}(y) &= E(X/Y)\end{aligned}$$

Thus, the minimum *mean-square error* estimation involves conditional expectation which is difficult to obtain numerically.

Let us consider a simpler version of the problem. We assume that  $\hat{X}(y) = ay$  and the estimation problem is to find the optimal value for  $a$ . Thus we have the *linear minimum mean-square error* criterion which minimizes  $E(X - aY)^2$ .

$$\begin{aligned}\frac{d}{da} E(X - aY)^2 &= 0 \\ \Rightarrow E \frac{d}{da} (X - aY)^2 &= 0 \\ \Rightarrow E(X - aY)Y &= 0 \\ \Rightarrow EeY &= 0\end{aligned}$$

where  $e$  is the estimation error.

The above result shows that for the *linear minimum mean-square error* criterion, estimation error is orthogonal to data. This result helps us in deriving optimal filters to estimate a random signal buried in noise.

The mean and variance also give some quantitative information about the bounds of RVs. Following inequalities are extremely useful in many practical problems.

### 1.20 Chebysev Inequality

Suppose  $X$  is a parameter of a manufactured item with known mean  $\mu_X$  and variance  $\sigma_X^2$ . The quality control department rejects the item if the absolute deviation of  $X$  from  $\mu_X$  is greater than  $2\sigma_X$ . What fraction of the manufacturing item does the quality control department reject? Can you roughly guess it?

The standard deviation gives us an intuitive idea how the random variable is distributed about the mean. This idea is more precisely expressed in the remarkable *Chebysev Inequality* stated below. For a random variable  $X$  with mean  $\mu_X$  and variance  $\sigma_X^2$

$$P\{|X - \mu_X| \geq \varepsilon\} \leq \frac{\sigma_X^2}{\varepsilon^2}$$

**Proof:**

$$\begin{aligned} \sigma_X^2 &= \int_{-\infty}^{\infty} (x - \mu_X)^2 f_X(x) dx \\ &\geq \int_{|X - \mu_X| \geq \varepsilon} (x - \mu_X)^2 f_X(x) dx \\ &\geq \int_{|X - \mu_X| \geq \varepsilon} \varepsilon^2 f_X(x) dx \\ &= \varepsilon^2 P\{|X - \mu_X| \geq \varepsilon\} \\ \therefore P\{|X - \mu_X| \geq \varepsilon\} &\leq \frac{\sigma_X^2}{\varepsilon^2} \end{aligned}$$

### 1.21 Markov Inequality

For a random variable  $X$  which take only nonnegative values

$$P\{X \geq a\} \leq \frac{E(X)}{a} \quad \text{where } a > 0.$$

$$\begin{aligned} E(X) &= \int_0^{\infty} x f_X(x) dx \\ &\geq \int_a^{\infty} x f_X(x) dx \\ &\geq \int_a^{\infty} a f_X(x) dx \\ &= a P\{X \geq a\} \end{aligned}$$

$$\therefore P\{X \geq a\} \leq \frac{E(X)}{a}$$

$$\textbf{Result: } P\{(X-k)^2 \geq a\} \leq \frac{E(X-k)^2}{a}$$

## 1.22 Convergence of a sequence of random variables

Let  $X_1, X_2, \dots, X_n$  be a sequence  $n$  independent and identically distributed random variables. Suppose we want to estimate the mean of the random variable on the basis of the observed data by means of the relation

$$\hat{\mu}_X = \frac{1}{n} \sum_{i=1}^n X_i$$

How closely does  $\hat{\mu}_X$  represent  $\mu_X$  as  $n$  is increased? How do we measure the closeness between  $\hat{\mu}_X$  and  $\mu_X$ ?

Notice that  $\hat{\mu}_X$  is a random variable. What do we mean by the statement  $\hat{\mu}_X$  converges to  $\mu_X$ ?

Consider a deterministic sequence  $x_1, x_2, \dots, x_n, \dots$ . The sequence converges to a limit  $x$  if correspond to any  $\varepsilon > 0$  we can find a positive integer  $m$  such that  $|x - x_n| < \varepsilon$  for  $n > m$ .

Convergence of a random sequence  $X_1, X_2, \dots, X_n, \dots$  cannot be defined as above.

A sequence of random variables is said to converge everywhere to  $X$  if

$$|X(\xi) - X_n(\xi)| \rightarrow 0 \text{ for } n > m \text{ and } \forall \xi.$$

## 1.23 Almost sure (a.s.) convergence or convergence with probability 1

For the random sequence  $X_1, X_2, \dots, X_n, \dots$

$\{X_n \rightarrow X\}$  this is an event.

$$\text{If } \begin{aligned} P\{s | X_n(s) \rightarrow X(s)\} &= 1 \quad \text{as } n \rightarrow \infty, \\ P\{s | |X_n(s) - X(s)| < \varepsilon \text{ for } n \geq m\} &= 1 \quad \text{as } m \rightarrow \infty, \end{aligned}$$

then the sequence is said to converge to  $X$  almost sure or with probability 1.

One important application is the **Strong Law of Large Numbers**:

If  $X_1, X_2, \dots, X_n, \dots$  are iid random variables, then  $\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu_X$  with probability 1 as  $n \rightarrow \infty$ .

## 1.24 Convergence in mean square sense

If  $E(X_n - X)^2 \rightarrow 0$  as  $n \rightarrow \infty$ , we say that the sequence converges to  $X$  in mean square (M.S).

### Example 7:

If  $X_1, X_2, \dots, X_n, \dots$  are iid random variables, then

$\frac{1}{n} \sum_{i=1}^N X_i \rightarrow \mu_X$  in the mean square as  $n \rightarrow \infty$ .

We have to show that  $\lim_{n \rightarrow \infty} E\left(\frac{1}{n} \sum_{i=1}^N X_i - \mu_X\right)^2 = 0$

Now,

$$\begin{aligned} E\left(\frac{1}{n} \sum_{i=1}^N X_i - \mu_X\right)^2 &= E\left(\frac{1}{n} \left(\sum_{i=1}^N (X_i - \mu_X)\right)\right)^2 \\ &= \frac{1}{n^2} \sum_{i=1}^N E(X_i - \mu_X)^2 + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1, j \neq i}^n E(X_i - \mu_X)(X_j - \mu_X) \\ &= \frac{n\sigma_X^2}{n^2} + 0 \text{ (Because of independence)} \\ &= \frac{\sigma_X^2}{n} \end{aligned}$$

$$\therefore \lim_{n \rightarrow \infty} E\left(\frac{1}{n} \sum_{i=1}^N X_i - \mu_X\right)^2 = 0$$

## 1.25 Convergence in probability

$P\{|X_n - X| > \varepsilon\}$  is a sequence of probability.  $X_n$  is said to convergent to  $X$  in probability if this sequence of probability is convergent that is

$$P\{|X_n - X| > \varepsilon\} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

If a sequence is convergent in mean, then it is convergent in probability also, because

$$P\{|X_n - X|^2 > \varepsilon^2\} \leq E(X_n - X)^2 / \varepsilon^2 \quad (\text{Markov Inequality})$$

We have

$$P\{|X_n - X| > \varepsilon\} \leq E(X_n - X)^2 / \varepsilon^2$$

If  $E(X_n - X)^2 \rightarrow 0$  as  $n \rightarrow \infty$ , (mean square convergent) then

$$P\{|X_n - X| > \varepsilon\} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

**Example 8:**

Suppose  $\{X_n\}$  be a sequence of random variables with

$$P(X_n = 1) = 1 - \frac{1}{n}$$

and

$$P(X_n = -1) = \frac{1}{n}$$

Clearly

$$P\{|X_n - 1| > \varepsilon\} = P\{X_n = -1\} = \frac{1}{n} \rightarrow 0$$

as  $n \rightarrow \infty$ .

Therefore  $\{X_n\} \xrightarrow{P} \{X = 0\}$

**1.26 Convergence in distribution**

The sequence  $X_1, X_2, \dots, X_n, \dots$  is said to converge to  $X$  in distribution if

$$F_{X_n}(x) \rightarrow F_X(x) \quad \text{as } n \rightarrow \infty.$$

Here the two distribution functions eventually coincide.

**1.27 Central Limit Theorem**

Consider independent and identically distributed random variables  $X_1, X_2, \dots, X_n$ .

Let  $Y = X_1 + X_2 + \dots + X_n$

Then  $\mu_Y = \mu_{X_1} + \mu_{X_2} + \dots + \mu_{X_n}$

And  $\sigma_Y^2 = \sigma_{X_1}^2 + \sigma_{X_2}^2 + \dots + \sigma_{X_n}^2$

The central limit theorem states that under very general conditions  $Y$  converges to  $N(\mu_Y, \sigma_Y^2)$  as  $n \rightarrow \infty$ . The conditions are:

1. The random variables  $X_1, X_2, \dots, X_n$  are independent with same mean and variance, but not identically distributed.
2. The random variables  $X_1, X_2, \dots, X_n$  are independent with different mean and same variance and not identically distributed.



## 1.28 Jointly Gaussian Random variables

Two random variables  $X$  and  $Y$  are called jointly Gaussian if their joint density function is

$$f_{X,Y}(x,y) = Ae^{-\frac{1}{2(1-\rho_{X,Y}^2)} \left[ \frac{(x-\mu_X)^2}{\sigma_X^2} - 2\rho_{XY} \frac{(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y} + \frac{(y-\mu_Y)^2}{\sigma_Y^2} \right]}$$

$$\text{where } A = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho_{X,Y}^2}}$$

### Properties:

(1) If  $X$  and  $Y$  are jointly Gaussian, then for any constants  $a$  and  $b$ , then the random variable

$Z$ , given by  $Z = aX + bY$  is Gaussian with mean  $\mu_Z = a\mu_X + b\mu_Y$  and variance

$$\sigma_Z^2 = a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\sigma_X\sigma_Y\rho_{X,Y}$$

(2) If two jointly Gaussian RVs are uncorrelated,  $\rho_{X,Y} = 0$  then they are statistically independent.

$$f_{X,Y}(x,y) = f_X(x)f_Y(y) \text{ in this case.}$$

(3) If  $f_{X,Y}(x,y)$  is a jointly Gaussian distribution, then the marginal densities

$$f_X(x) \text{ and } f_Y(y) \text{ are also Gaussian.}$$

(4) If  $X$  and  $Y$  are joint by Gaussian random variables then the optimum nonlinear estimator  $\hat{X}$  of  $X$  that minimizes the mean square error  $\xi = E\{[X - \hat{X}]^2\}$  is a linear estimator  $\hat{X} = aY$

## CHAPTER - 2 : REVIEW OF RANDOM PROCESS

### 2.1 Introduction

Recall that a random variable maps each sample point in the sample space to a point in the real line. A random process maps each sample point to a waveform.

- A random process can be defined as an indexed family of random variables  $\{X(t), t \in T\}$  where  $T$  is an index set which may be discrete or continuous usually denoting time.
- The random process is defined on a common probability space  $\{S, \mathfrak{F}, P\}$ .
- A random process is a function of the sample point  $\xi$  and index variable  $t$  and may be written as  $X(t, \xi)$ .
- For a fixed  $t (= t_0)$ ,  $X(t_0, \xi)$  is a random variable.
- For a fixed  $\xi (= \xi_0)$ ,  $X(t, \xi_0)$  is a single realization of the random process and is a deterministic function.
- When both  $t$  and  $\xi$  are varying we have the random process  $X(t, \xi)$ .

The random process  $X(t, \xi)$  is normally denoted by  $X(t)$ .

We can define a discrete random process  $X[n]$  on discrete points of time. Such a random process is more important in practical implementations.

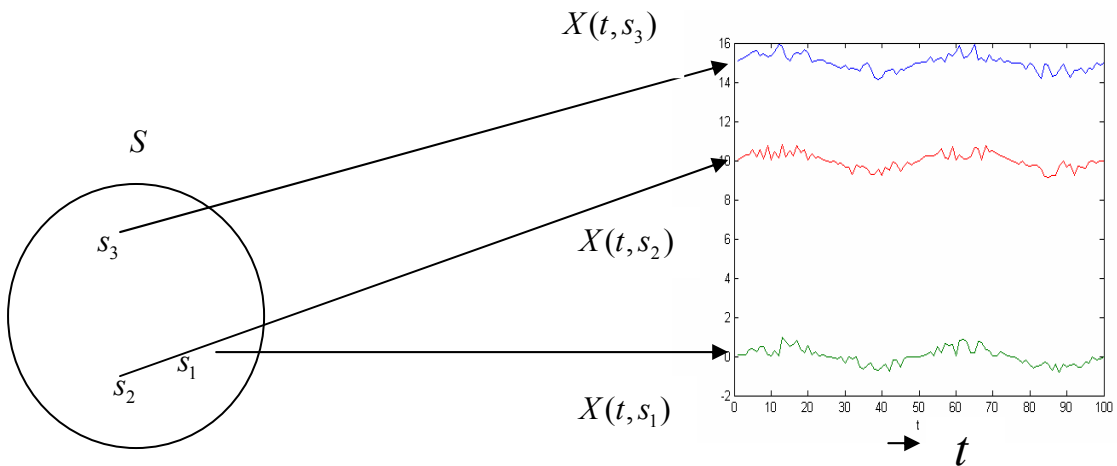


Figure Random Process

## 2.2 How to describe a random process?

To describe  $X(t)$  we have to use joint density function of the random variables at different  $t$ .

For any positive integer  $n$ ,  $X(t_1), X(t_2), \dots, X(t_n)$  represents  $n$  jointly distributed random variables. Thus a random process can be described by the joint distribution function  $F_{X(t_1), X(t_2), \dots, X(t_n)}(x_1, x_2, \dots, x_n) = F(x_1, x_2, \dots, x_n, t_1, t_2, \dots, t_n), \forall n \in N$  and  $\forall t_n \in T$

Otherwise we can determine all the possible moments of the process.

$E(X(t)) = \mu_x(t)$  = mean of the random process at  $t$ .

$R_X(t_1, t_2) = E(X(t_1)X(t_2))$  = autocorrelation function at  $t_1, t_2$

$R_X(t_1, t_2, t_3) = E(X(t_1), X(t_2), X(t_3))$  = Triple correlation function at  $t_1, t_2, t_3$ , etc.

We can also define the auto-covariance function  $C_X(t_1, t_2)$  of  $X(t)$  given by

$$\begin{aligned} C_X(t_1, t_2) &= E(X(t_1) - \mu_X(t_1))(X(t_2) - \mu_X(t_2)) \\ &= R_X(t_1, t_2) - \mu_X(t_1)\mu_X(t_2) \end{aligned}$$

### Example 1:

#### (a) Gaussian Random Process

For any positive integer  $n$ ,  $X(t_1), X(t_2), \dots, X(t_n)$  represent  $n$  jointly random variables. These  $n$  random variables define a random vector  $\mathbf{X} = [X(t_1), X(t_2), \dots, X(t_n)]'$ . The process  $X(t)$  is called Gaussian if the random vector  $[X(t_1), X(t_2), \dots, X(t_n)]'$  is jointly Gaussian with the joint density function given by

$$f_{X(t_1), X(t_2), \dots, X(t_n)}(x_1, x_2, \dots, x_n) = \frac{e^{-\frac{1}{2}\mathbf{x}'\mathbf{C}_X^{-1}\mathbf{x}}}{\left(\sqrt{2\pi}\right)^n \sqrt{\det(\mathbf{C}_X)}} \text{ where } \mathbf{C}_X = E(\mathbf{X} - \boldsymbol{\mu}_X)(\mathbf{X} - \boldsymbol{\mu}_X)'$$

$$\text{and } \boldsymbol{\mu}_X = E(\mathbf{X}) = [E(X_1), E(X_2), \dots, E(X_n)]'.$$

#### (b) Bernouli Random Process

#### (c) A sinusoid with a random phase.

## 2.3 Stationary Random Process

A random process  $X(t)$  is called strict-sense stationary if its probability structure is invariant with time. In terms of the joint distribution function

$$F_{X(t_1), X(t_2), \dots, X(t_n)}(x_1, x_2, \dots, x_n) = F_{X(t_1+t_0), X(t_2+t_0), \dots, X(t_n+t_0)}(x_1, x_2, \dots, x_n) \quad \forall n \in N \text{ and } \forall t_0, t_n \in T$$

For  $n = 1$ ,

$$F_{X(t_1)}(x_1) = F_{X(t_1+t_0)}(x_1) \quad \forall t_0 \in T$$

Let us assume  $t_0 = -t_1$

$$\begin{aligned} F_{X(t_1)}(x_1) &= F_{X(0)}(x_1) \\ \Rightarrow EX(t_1) &= EX(0) = \mu_X(0) = \text{constant} \end{aligned}$$

For  $n = 2$ ,

$$F_{X(t_1), X(t_2)}(x_1, x_2) = F_{X(t_1+t_0), X(t_2+t_0)}(x_1, x_2)$$

Put  $t_0 = -t_2$

$$\begin{aligned} F_{X(t_1), X(t_2)}(x_1, x_2) &= F_{X(t_1-t_2), X(0)}(x_1, x_2) \\ \Rightarrow R_X(t_1, t_2) &= R_X(t_1 - t_2) \end{aligned}$$

A random process  $X(t)$  is called **wide sense stationary process (WSS)** if

$$\begin{aligned} \mu_X(t) &= \text{constant} \\ R_X(t_1, t_2) &= R_X(t_1 - t_2) \text{ is a function of time lag.} \end{aligned}$$

For a Gaussian random process, WSS implies strict sense stationarity, because this process is completely described by the mean and the autocorrelation functions.

The autocorrelation function  $R_X(\tau) = EX(t+\tau)X(t)$  is a crucial quantity for a WSS process.

- $R_X(0) = EX^2(t)$  is the mean-square value of the process.
- $R_X(-\tau) = R_X(\tau)$  for real process (for a complex process  $X(t)$ ,  $R_X(-\tau) = R_X^*(\tau)$ )
- $|R_X(\tau)| \leq R_X(0)$  which follows from the Schwartz inequality

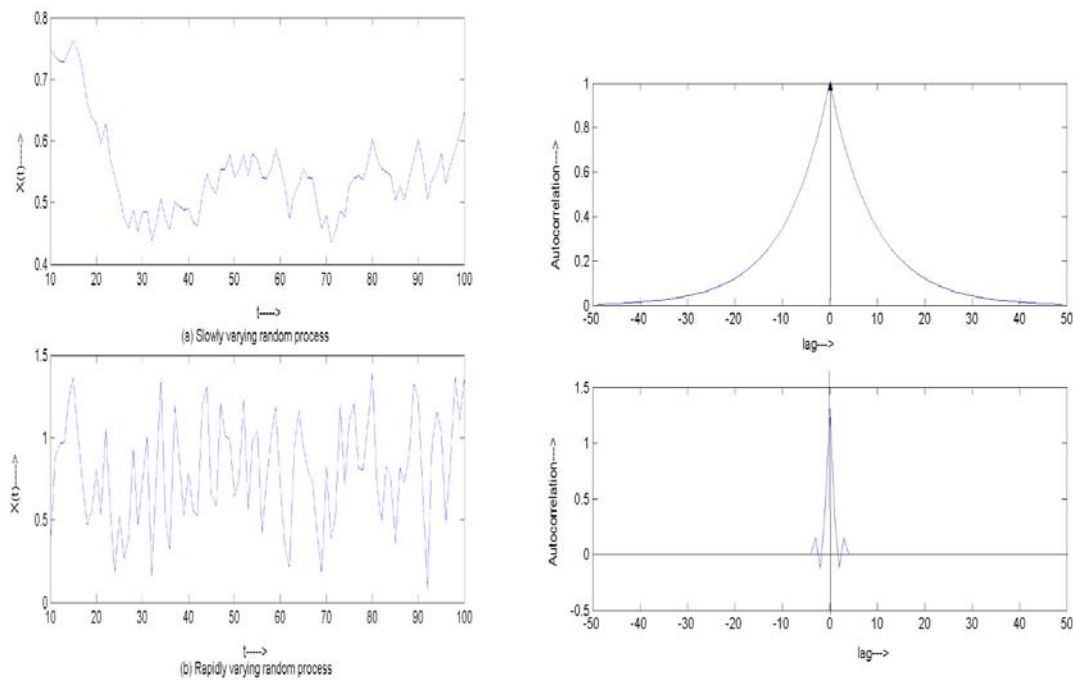
$$\begin{aligned} R_X^2(\tau) &= \{EX(t)X(t+\tau)\}^2 \\ &= |\langle X(t), X(t+\tau) \rangle|^2 \\ &\leq \|X(t)\|^2 \|X(t+\tau)\|^2 \\ &= EX^2(t)EX^2(t+\tau) \\ &= R_X^2(0)R_X^2(0) \\ \therefore |R_X(\tau)| &\leq R_X(0) \end{aligned}$$

- $R_X(\tau)$  is a positive semi-definite function in the sense that for any positive integer  $n$  and real  $a_j$ ,  $\sum_{i=1}^n \sum_{j=1}^n a_i a_j R_X(t_i, t_j) \geq 0$
- If  $X(t)$  is periodic (in the mean square sense or any other sense like with probability 1), then  $R_X(\tau)$  is also periodic.

For a discrete random sequence, we can define the autocorrelation sequence similarly.

- If  $R_X(\tau)$  drops quickly, then the signal samples are less correlated which in turn means that the signal has lot of changes with respect to time. Such a signal has high frequency components. If  $R_X(\tau)$  drops slowly, the signal samples are highly correlated and such a signal has less high frequency components.
- $R_X(\tau)$  is directly related to the frequency domain representation of WSS process.

The following figure illustrates the above concepts



**Figure** Frequency Interpretation of Random process: for slowly varying random process Autocorrelation decays slowly

## 2.4 Spectral Representation of a Random Process

*How to have the frequency-domain representation of a random process?*

- Wiener (1930) and Khinchin (1934) independently discovered the spectral representation of a random process. Einstein (1914) also used the concept.
- Autocorrelation function and power spectral density forms a Fourier transform pair

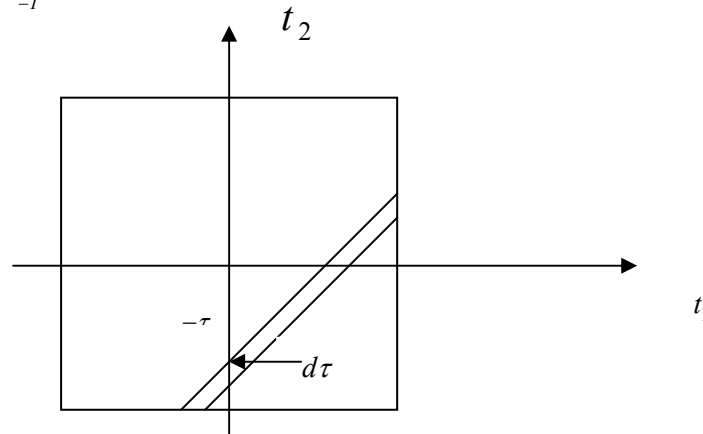
Lets define

$$X_T(t) = X(t) \quad -T < t < T$$

$$= 0 \quad \text{otherwise}$$

as  $t \rightarrow \infty$ ,  $X_T(t)$  will represent the random process  $X(t)$ .

Define  $X_T(\omega) = \int_{-T}^T X_T(t) e^{-j\omega t} dt$  in mean square sense.



$$E \frac{X_T(\omega) X_T^*(\omega)}{2T} = E \frac{|X_T(\omega)|^2}{2T} = \frac{1}{2T} \int_{-T}^T \int_{-T}^T E X_T(t_1) X_T(t_2) e^{-j\omega t_1} e^{+j\omega t_2} dt_1 dt_2$$

$$= \frac{1}{2T} \int_{-T}^T \int_{-T}^T R_X(t_1 - t_2) e^{-j\omega(t_1 - t_2)} dt_1 dt_2$$

$$= \frac{1}{2T} \int_{-2T}^{2T} R_X(\tau) e^{-j\omega \tau} (2T - |\tau|) d\tau$$

Substituting  $t_1 - t_2 = \tau$  so that  $t_2 = t_1 - \tau$  is a line, we get

$$E \frac{X_T(\omega) X_T^*(\omega)}{2T} = \int_{-2T}^{2T} R_X(\tau) e^{-j\omega \tau} \left(1 - \frac{|\tau|}{2T}\right) d\tau$$

If  $R_X(\tau)$  is integrable then as  $T \rightarrow \infty$ ,

$$\lim_{T \rightarrow \infty} \frac{E |X_T(\omega)|^2}{2T} = \int_{-\infty}^{\infty} R_X(\tau) e^{-j\omega \tau} d\tau$$

$\frac{E|X_T(\omega)|^2}{2T}$  = contribution to average power at freq  $\omega$  and is called the power spectral density.

Thus 
$$S_X(\omega) = \int_{-\infty}^{\infty} R_X(\tau) e^{-j\omega\tau} d\tau$$

and 
$$R_X(\tau) = \int_{-\infty}^{\infty} S_X(\omega) e^{j\omega\tau} d\omega$$

### **Properties**

- $EX^2(t) = R_X(0) = \int_{-\infty}^{\infty} S_X(\omega) d\omega$  = average power of the process.
- The average power in the band  $(w_1, w_2)$  is  $\int_{w_1}^{w_2} S_X(w) dw$
- $R_X(\tau)$  is real and even  $\Rightarrow S_X(\omega)$  is real, even.
- From the definition  $S_X(w) = \lim_{T \rightarrow \infty} \frac{E|X_T(\omega)|^2}{2T}$  is always positive.
- $h_X(w) = \frac{S_X(\omega)}{EX^2(t)}$  = normalised power spectral density and has properties of PDF, (always +ve and area=1).

## **2.5 Cross-correlation & Cross power Spectral Density**

Consider two real random processes  $X(t)$  and  $Y(t)$ .

Joint stationarity of  $X(t)$  and  $Y(t)$  implies that the joint densities are invariant with shift of time.

The **cross-correlation function**  $R_{XY}(\tau)$  for a jointly wss processes  $X(t)$  and  $Y(t)$  is defined as

$$R_{X,Y}(\tau) = E X(t + \tau)Y(t)$$

so that  $R_{YX}(\tau) = E Y(t + \tau)X(t)$

$$= E X(t)Y(t + \tau)$$

$$= R_{X,Y}(-\tau)$$

$$\therefore R_{YX}(\tau) = R_{X,Y}(-\tau)$$

Cross power spectral density

$$S_{X,Y}(w) = \int_{-\infty}^{\infty} R_{X,Y}(\tau) e^{-jw\tau} d\tau$$

For real processes  $X(t)$  and  $Y(t)$

$$S_{X,Y}(w) = S_{Y,X}^*(w)$$

The Wiener-Khinchin theorem is also valid for discrete-time random processes.

If we define  $R_X[m] = E X[n+m] X[n]$

Then corresponding PSD is given by

$$S_X(w) = \sum_{m=-\infty}^{\infty} R_X[m] e^{-j\omega m} \quad -\pi \leq w \leq \pi$$

$$\text{or } S_X(f) = \sum_{m=-\infty}^{\infty} R_X[m] e^{-j2\pi f m} \quad -1 \leq f \leq 1$$

$$\therefore R_X[m] = \frac{1}{2\pi} \int_{-\pi}^{\pi} S_X(w) e^{j\omega m} dw$$

For a discrete sequence the generalized PSD is defined in the  $z$ -domain as follows

$$S_X(z) = \sum_{m=-\infty}^{\infty} R_X[m] z^{-m}$$

If we sample a stationary random process uniformly we get a stationary random sequence.

Sampling theorem is valid in terms of PSD.

### **Examples 2:**

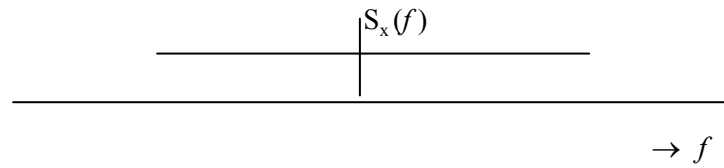
$$(1) R_X(\tau) = e^{-a|\tau|} \quad a > 0$$

$$S_X(w) = \frac{2a}{a^2 + w^2} \quad -\infty < w < \infty$$

$$(2) R_X(m) = a^{|m|} \quad |a| > 0$$

$$S_X(w) = \frac{1 - a^2}{1 - 2a \cos w + a^2} \quad -\pi \leq w \leq \pi$$

## **2.6 White noise process**



A white noise process  $X(t)$  is defined by

$$S_X(f) = \frac{N}{2} \quad -\infty < f < \infty$$

The corresponding autocorrelation function is given by

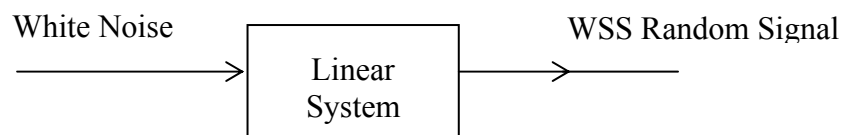
$$R_X(\tau) = \frac{N}{2} \delta(\tau) \quad \text{where } \delta(\tau) \text{ is the Dirac delta.}$$

The average power of white noise

$$P_{avg} = \int_{-\infty}^{\infty} \frac{N}{2} df \rightarrow \infty$$



- Samples of a white noise process are uncorrelated.
- White noise is an mathematical abstraction, it cannot be realized since it has infinite power
- If the system band-width(BW) is sufficiently narrower than the noise BW and noise PSD is flat , we can model it as a white noise process. Thermal and shot noise are well modelled as white Gaussian noise, since they have very flat psd over very wide band (GHzs
- For a zero-mean white noise process, the correlation of the process at any lag  $\tau \neq 0$  is zero.
- White noise plays a key role in random signal modelling.
- Similar role as that of the impulse function in the modeling of deterministic signals.



## 2.7 White Noise Sequence

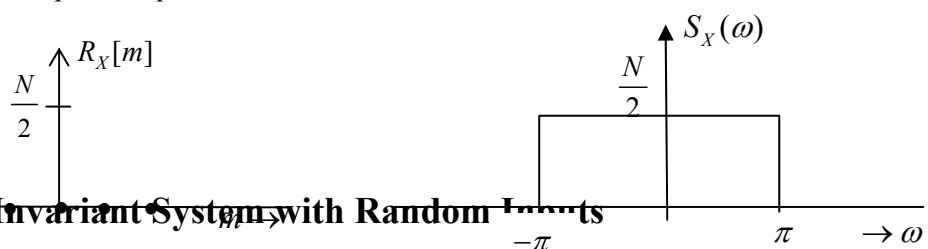
For a white noise sequence  $x[n]$ ,

$$S_X(\omega) = \frac{N}{2} \quad -\pi \leq \omega \leq \pi$$

Therefore

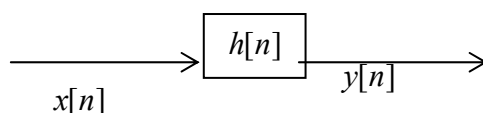
$$R_X(m) = \frac{N}{2} \delta(m)$$

where  $\delta(m)$  is the unit impulse sequence.



## 2.8 Linear Shift Invariant System with Random Inputs

Consider a discrete-time linear system with impulse response  $h[n]$ .



$$y[n] = x[n] * h[n]$$

$$E y[n] = E x[n] * h[n]$$

For stationary input  $x[n]$

$$\mu_Y = E y[n] = \mu_X * h[n] = \mu_X \sum_{k=0}^l h[k]$$

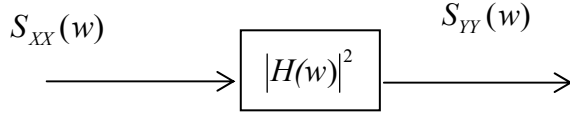
where  $l$  is the length of the impulse response sequence

$$\begin{aligned} R_Y[m] &= E y[n]y[n-m] \\ &= E(x[n] * h[n]) * (x[n-m] * h[n-m]) \\ &= R_X[m] * h[m] * h[-m] \end{aligned}$$

$R_Y[m]$  is a function of lag  $m$  only.

From above we get

$$S_Y(w) = |H(w)|^2 S_X(w)$$



### **Example 3:**

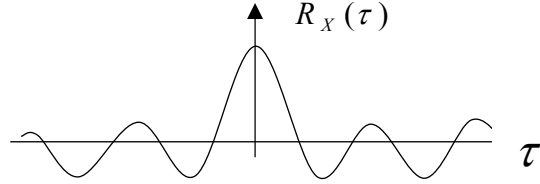
Suppose

$$\begin{aligned} H(w) &= 1 \quad -w_c \leq w \leq w_c \\ &= 0 \quad \text{otherwise} \end{aligned}$$

$$S_X(w) = \frac{N}{2} \quad -\infty \leq w \leq \infty$$

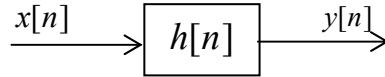
Then  $S_Y(w) = \frac{N}{2} \quad -w_c \leq w \leq w_c$

and  $R_Y(\tau) = \frac{N}{2} \text{sinc}(w_c \tau)$



- Note that though the input is an uncorrelated process, the output is a correlated process.

Consider the case of the discrete-time system with a random sequence  $x[n]$  as an input.



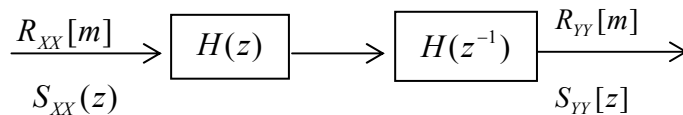
$$R_Y[m] = R_X[m] * h[m] * h[-m]$$

Taking the  $z$  – transform, we get

$$S_Y(z) = S_X(z)H(z)H(z^{-1})$$

Notice that if  $H(z)$  is causal, then  $H(z^{-1})$  is anti causal.

Similarly if  $H(z)$  is minimum-phase then  $H(z^{-1})$  is maximum-phase.



**Example 4:**

If  $H(z) = \frac{1}{1-\alpha z^{-1}}$  and  $x[n]$  is a unity-variance white-noise sequence, then

$$\begin{aligned} S_{YY}(z) &= H(z)H(z^{-1}) \\ &= \left(\frac{1}{1-\alpha z^{-1}}\right)\left(\frac{1}{1-\alpha z}\right)\frac{1}{2\pi} \end{aligned}$$

By partial fraction expansion and inverse  $z$  – transform, we get

$$R_Y[m] = \frac{1}{1-\alpha^2} \alpha^{|m|}$$

**2.9 Spectral factorization theorem**

A stationary random signal  $X[n]$  that satisfies the Paley Wiener condition

$\int_{-\pi}^{\pi} |\ln S_X(w)| dw < \infty$  can be considered as an output of a linear filter fed by a white noise sequence.

If  $S_X(w)$  is an analytic function of  $w$ ,

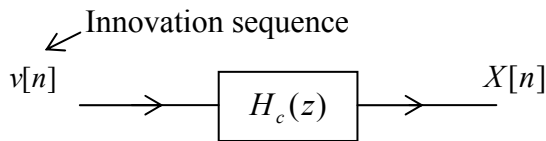
and  $\int_{-\pi}^{\pi} |\ln S_X(w)| dw < \infty$ , then  $S_X(z) = \sigma_v^2 H_c(z) H_a(z)$

where

$H_c(z)$  is the causal minimum phase transfer function

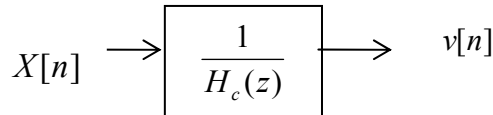
$H_a(z)$  is the anti-causal maximum phase transfer function

and  $\sigma_v^2$  a constant and interpreted as the variance of a white-noise sequence.



**Figure** Innovation Filter

Minimum phase filter => the corresponding inverse filter exists.



**Figure** whitening filter

Since  $\ln S_{XX}(z)$  is analytic in an annular region  $\rho < |z| < \frac{1}{\rho}$ ,

$$\ln S_{XX}(z) = \sum_{k=-\infty}^{\infty} c[k] z^{-k}$$

where  $c[k] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \ln S_{XX}(w) e^{iwn} dw$  is the  $k$ th order cepstral coefficient.

For a real signal  $c[k] = c[-k]$

and  $c[0] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \ln S_{XX}(w) dw$

$$\begin{aligned} S_{XX}(z) &= e^{\sum_{k=-\infty}^{\infty} c[k]z^{-k}} \\ &= e^{c[0]} e^{\sum_{k=1}^{\infty} c[k]z^{-k}} e^{\sum_{k=-\infty}^{-1} c[k]z^{-k}} \end{aligned}$$

$$\begin{aligned} \text{Let } H_C(z) &= e^{\sum_{k=1}^{\infty} c[k]z^{-k}} \quad |z| > \rho \\ &= 1 + h_c(1)z^{-1} + h_c(2)z^{-2} + \dots \end{aligned}$$

$$(\because h_c[0] = \lim_{z \rightarrow \infty} H_C(z) = 1)$$

$H_C(z)$  and  $\ln H_C(z)$  are both analytic

$\Rightarrow H_C(z)$  is a **minimum phase filter**.

Similarly let

$$\begin{aligned} H_a(z) &= e^{\sum_{k=-\infty}^{-1} c(k)z^{-k}} \\ &= e^{\sum_{k=1}^{\infty} c(k)z^k} = H_C(z^{-1}) \quad |z| < \frac{1}{\rho} \end{aligned}$$

Therefore,

$$S_{XX}(z) = \sigma_v^2 H_C(z) H_C(z^{-1})$$

where  $\sigma_v^2 = e^{c(0)}$

### Salient points

- $S_{XX}(z)$  can be factorized into a minimum-phase and a maximum-phase factors i.e.  $H_C(z)$  and  $H_C(z^{-1})$ .
- In general spectral factorization is difficult, however for a signal with rational power spectrum, spectral factorization can be easily done.
- Since is a minimum phase filter,  $\frac{1}{H_C(z)}$  exists ( $\Rightarrow$  stable), therefore we can have a filter  $\frac{1}{H_C(z)}$  to filter the given signal to get the innovation sequence.
- $X[n]$  and  $v[n]$  are related through an invertible transform; so they contain the same information.

## 2.10 Wold's Decomposition

Any WSS signal  $X[n]$  can be decomposed as a sum of two mutually orthogonal processes

- a regular process  $X_r[n]$  and a predictable process  $X_p[n]$ ,  $X[n] = X_r[n] + X_p[n]$
- $X_r[n]$  can be expressed as the output of linear filter using a white noise sequence as input.
- $X_p[n]$  is a predictable process, that is, the process can be predicted from its own past with zero prediction error.

## CHAPTER - 3: RANDOM SIGNAL MODELLING

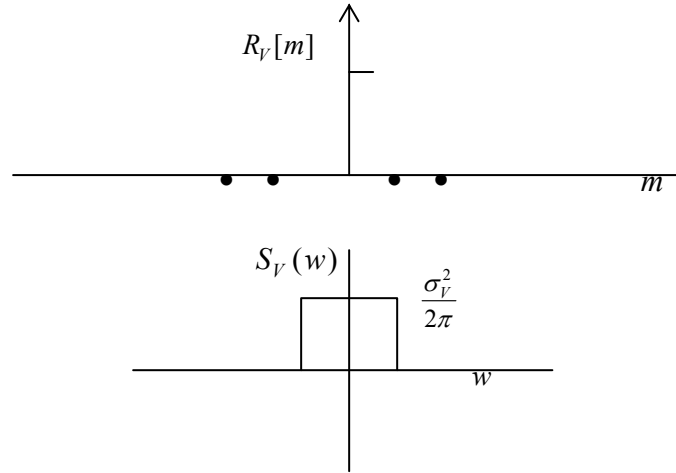
### 3.1 Introduction

The spectral factorization theorem enables us to model a regular random process as an output of a linear filter with white noise as input. Different models are developed using different forms of linear filters.

- These models are mathematically described by linear constant coefficient difference equations.
- In statistics, random-process modeling using difference equations is known as *time series analysis*.

### 3.2 White Noise Sequence

The simplest model is the white noise  $v[n]$ . We shall assume that  $v[n]$  is of 0-mean and variance  $\sigma_v^2$ .



### 3.3 Moving Average model $MA(q)$ model



The difference equation model is

$$X[n] = \sum_{i=0}^q b_i v[n-i]$$

$$\mu_e = 0 \Rightarrow \mu_X = 0$$

and  $v[n]$  is an uncorrelated sequence means

$$\sigma_X^2 = \sum_{i=0}^q b_i^2 \sigma_v^2$$

The autocorrelations are given by

$$\begin{aligned}
R_X[m] &= E X[n] X[n-m] \\
&= \sum_{i=0}^q \sum_{j=0}^q b_i b_j E v[n-i] v[n-m-j] \\
&= \sum_{i=0}^q \sum_{j=0}^q b_i b_j R_V[m-i+j]
\end{aligned}$$

Noting that  $R_V[m] = \sigma_V^2 \delta[m]$ , we get

$$\begin{aligned}
R_V[m] &= \sigma_V^2 \text{ when} \\
m-i+j &= 0 \\
\Rightarrow i &= m+j
\end{aligned}$$

The maximum value for  $m+j$  is  $q$  so that

$$R_X[m] = \sum_{j=0}^{q-m} b_j b_{j+m} \sigma_V^2 \quad 0 \leq m \leq q$$

and

$$R_X[-m] = R_X[m]$$

Writing the above two relations together 
$$R_X[m] = \sum_{j=0}^{q-|m|} b_j b_{j+|m|} \sigma_V^2 \quad |m| \leq q$$
  

$$= 0 \text{ otherwise}$$

Notice that,  $R_X[m]$  is related by a nonlinear relationship with model parameters. Thus finding the model parameters is not simple.

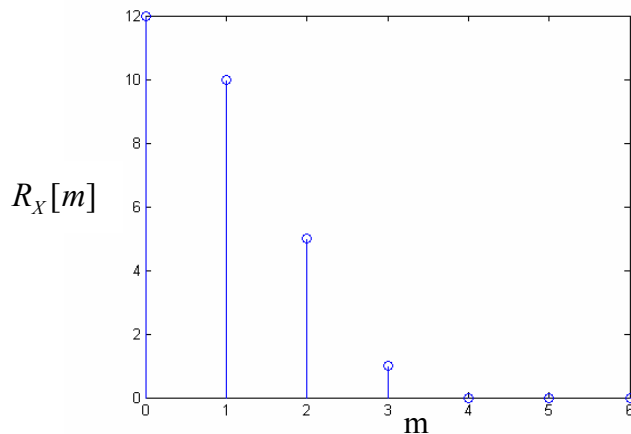
The power spectral density is given by

$$S_X(w) = \frac{\sigma_V^2}{2\pi} |B(w)|^2, \text{ where } B(w) = b_0 + b_1 e^{-jw} + \dots b_q e^{-jqw}$$

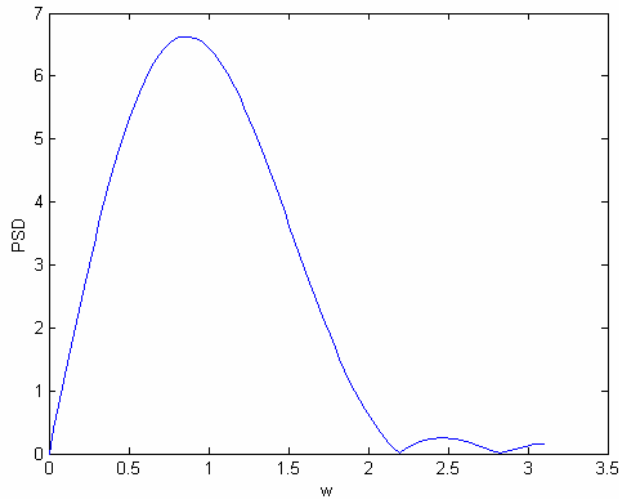
FIR system will give some zeros. So if the spectrum has some valleys then MA will fit well.

### 3.3.1 Test for MA process

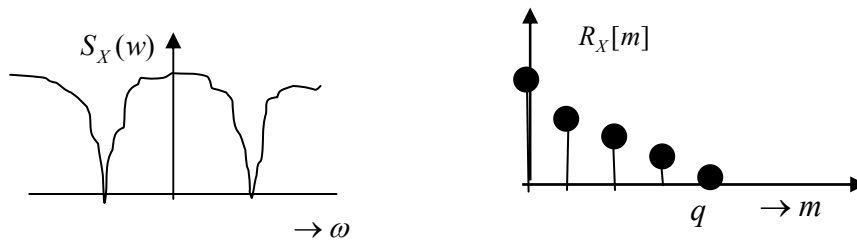
$R_X[m]$  becomes zero suddenly after some value of  $m$ .



**Figure:** Autocorrelation function of a MA process



**Figure:** Power spectrum of a MA process



### **Example 1: MA(1) process**

$$X[n] = b_1 v[n-1] + b_0 v[n]$$

Here the parameters  $b_0$  and  $b_1$  are to be determined.

We have

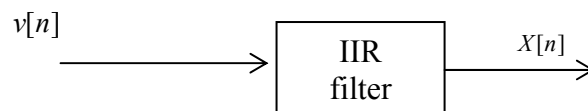
$$\sigma_X^2 = b_1^2 + b_0^2$$

$$R_X[1] = b_1 b_0$$

From above  $b_0$  and  $b_1$  can be calculated using the variance and autocorrelation at lag 1 of the signal.

## **3.4 Autoregressive Model**

In time series analysis it is called AR(p) model.



The model is given by the difference equation

$$X[n] = \sum_{i=1}^p a_i X[n-i] + v[n]$$

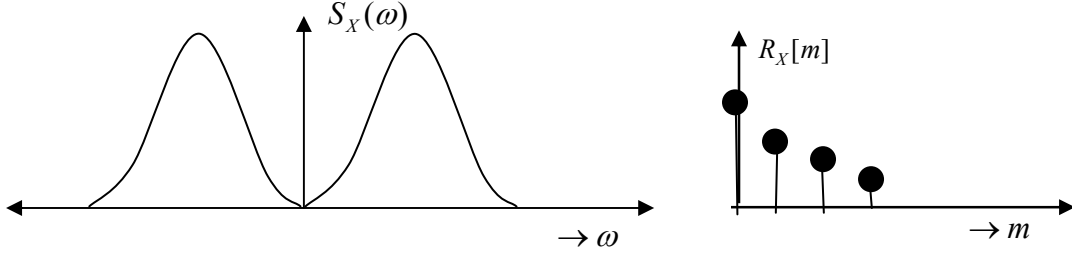
The transfer function  $A(w)$  is given by



$$A(w) = \frac{1}{1 - \sum_{i=1}^n a_i e^{-j\omega i}}$$

with  $a_0 = 1$  (all poles model) and  $S_X(w) = \frac{\sigma_e^2}{2\pi |A(w)|^2}$

If there are sharp peaks in the spectrum, the AR(p) model may be suitable.



The autocorrelation function  $R_X[m]$  is given by

$$\begin{aligned} R_X[m] &= E X[n] X[n-m] \\ &= \sum_{i=1}^p a_i E X[n-i] X[n-m] + E v[n] X[n-m] \\ &= \sum_{i=1}^p a_i R_X[m-i] + \sigma_v^2 \delta[m] \end{aligned}$$

$$\therefore R_X[m] = \sum_{i=1}^p a_i R_X[m-i] + \sigma_v^2 \delta[m] \quad \forall m \in I$$

The above relation gives a set of linear equations which can be solved to find  $a_i$ s.

These sets of equations are known as Yule-Walker Equation.

### **Example 2: AR(1) process**

$$X[n] = a_1 X[n-1] + v[n]$$

$$R_X[m] = a_1 R_X[m-1] + \sigma_v^2 \delta[m]$$

$$\therefore R_X[0] = a_1 R_X[-1] + \sigma_v^2 \quad (1)$$

$$\text{and } R_X[1] = a_1 R_X[0]$$

$$\text{so that } a_1 = \frac{R_X[1]}{R_X[0]}$$

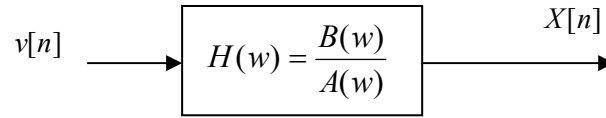
$$\text{From (1) } \sigma_v^2 = R_X[0] = \frac{\sigma_v^2}{1-a^2}$$

After some arithmetic we get

$$R_X[m] = \frac{a^{|m|} \sigma_v^2}{1-a^2}$$

### 3.5 ARMA(p,q) – Autoregressive Moving Average Model

Under the most practical situation, the process may be considered as an output of a filter that has both zeros and poles.



The model is given by

$$x[n] = \sum_{i=1}^p a_i x[n-i] + \sum_{i=0}^q b_i v[n-i] \quad (\text{ARMA 1})$$

and is called the  $ARMA(p, q)$  model.

The transfer function of the filter is given by

$$H(\omega) = \frac{B(\omega)}{A(\omega)}$$

$$S_X(\omega) = \frac{|B(\omega)|^2 \sigma_v^2}{|A(\omega)|^2 2\pi}$$

How do get the model parameters?

For  $m \geq \max(p, q + 1)$ , there will be no contributions from  $b_i$  terms to  $R_X[m]$ .

$$R_X[m] = \sum_{i=1}^p a_i R_X[m-i] \quad m \geq \max(p, q + 1)$$

From a set of  $p$  Yule Walker equations,  $a_i$  parameters can be found out.

Then we can rewrite the equation

$$\tilde{X}[n] = X[n] + \sum_{i=1}^p a_i X[n-i]$$

$$\therefore \tilde{X}[n] = \sum_{i=0}^q b_i v[n-i]$$

From the above equation  $b_i$ s can be found out.

The  $ARMA(p, q)$  is an economical model. Only  $AR(p)$  only  $MA(q)$  model may require a large number of model parameters to represent the process adequately. This concept in model building is known as ***the parsimony of parameters***.

The difference equation of the  $ARMA(p, q)$  model, given by eq. (ARMA 1) can be reduced to  $p$  first-order difference equation give a state space representation of the random process as follows:

$$\mathbf{z}[n] = \mathbf{A}\mathbf{z}[n-1] + \mathbf{B}u[n]$$

$$X[n] = \mathbf{C}\mathbf{z}[n]$$

where

$$\mathbf{z}[n] = [x[n] \ X[n-1] \dots X[n-p]]'$$

$$\mathbf{A} = \begin{bmatrix} a_1 & a_2 & \dots & a_p \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{bmatrix}, \quad \mathbf{B} = [1 \ 0 \dots 0]' \text{ and}$$

$$\mathbf{C} = [b_0 \ b_1 \dots b_q]$$

Such representation is convenient for analysis.

### 3.6 General $ARMA(p, q)$ Model building Steps

- Identification of p and q.
  - Estimation of model parameters.
  - Check the modeling error.
  - If it is white noise then stop.
- Else select new values for p and q  
and repeat the process.

### 3.7 Other model: To model nonstationary random processes

- **ARIMA model:** Here after differencing the data can be fed to an ARMA model.
- **SARMA model:** Seasonal ARMA model etc. Here the signal contains a seasonal fluctuation term. The signal after differencing by step equal to the seasonal period becomes stationary and ARMA model can be fitted to the resulting data.

## **SECTION – II**

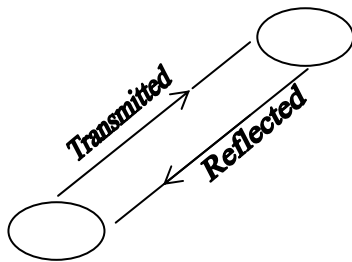
### **ESTIMATION THEORY**

## CHAPTER – 4: ESTIMATION THEORY

### 4.1 Introduction

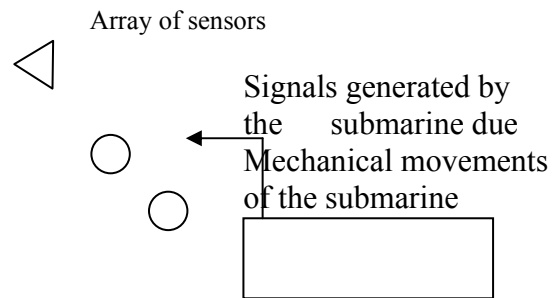
- For speech, we have LPC (linear predictive code) model, the LPC-parameters are to be estimated from observed data.
- We may have to estimate the correct value of a signal from the noisy observation.

In RADAR signal processing



- Estimate the target, target distance from the observed data

In sonar signal processing



- Estimate the location of the submarine.

Generally *estimation* includes parameter *estimation* and signal *estimation*.

We will discuss the problem of parameter estimation here.

We have a sequence of observable random variables  $X_1, X_2, \dots, X_N$ , represented by the vector

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$$

$\mathbf{X}$  is governed by a joint density junction which depends on some unobservable parameter  $\theta$  given by

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_N | \theta) = f_{\mathbf{X}}(\mathbf{x} | \theta)$$

where  $\theta$  may be deterministic or random. Our aim is to make an inference on  $\theta$  from an observed sample of  $X_1, X_2, \dots, X_N$ .

An estimator  $\hat{\theta}(\mathbf{X})$  is a rule by which we guess about the value of an unknown  $\theta$  on the basis of  $\mathbf{X}$ .

$\hat{\theta}(\mathbf{X})$  is a random, being a function of random variables.

For a particular observation  $x_1, x_2, \dots, x_N$ , we get what is known as an estimate (not estimator)

Let  $X_1, X_2, \dots, X_N$  be a sequence of independent and identically distributed (iid) random variables with mean  $\mu_X$  and variance  $\sigma_X^2$ .

$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$  is an estimator for  $\mu_X$ .

$\hat{\sigma}_X^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu}_X)^2$  is an estimator for  $\sigma_X^2$ .

An estimator is a function of the random sequence  $X_1, X_2, \dots, X_N$  and if it does not involve any unknown parameters. Such a function is generally called a statistic.

## 4.2 Properties of the Estimator

A good estimator should satisfy some properties. These properties are described in terms of the mean and variance of the estimator.

## 4.3 Unbiased estimator

An estimator  $\hat{\theta}$  of  $\theta$  is said to be unbiased if and only if  $E\hat{\theta} = \theta$ .

The quantity  $E\hat{\theta} - \theta$  is called the bias of the estimator.

Unbiased ness is necessary but not sufficient to make an estimator a good one.

Consider  $\hat{\sigma}_1^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu}_X)^2$

and  $\hat{\sigma}_2^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \hat{\mu}_X)^2$

for an iid random sequence  $X_1, X_2, \dots, X_N$ .

We can show that  $\hat{\sigma}_2^2$  is an unbiased estimator.

$$\begin{aligned} E \sum_{i=1}^N (X_i - \hat{\mu}_X)^2 &= E \sum_{i=1}^N (X_i - \mu_X + \mu_X - \hat{\mu}_X)^2 \\ &= E \sum_{i=1}^N \{ (X_i - \mu_X)^2 + (\mu_X - \hat{\mu}_X)^2 + 2(X_i - \mu_X)(\mu_X - \hat{\mu}_X) \} \end{aligned}$$

Now  $E(X_i - \mu_X)^2 = \sigma^2$

$$\begin{aligned}
\text{and } E(\mu_X - \hat{\mu}_X)^2 &= E\left(\mu_X - \frac{\sum X_i}{N}\right)^2 \\
&= \frac{E}{N^2} (N\mu_X - \sum X_i)^2 \\
&= \frac{E}{N^2} (\sum (X_i - \mu_X))^2 \\
&= \frac{E}{N^2} \sum (X_i - \mu_X)^2 + \sum_i \sum_{j \neq i} E(X_i - \mu_X)(X_j - \mu_X) \\
&= \frac{E}{N^2} \sum (X_i - \mu_X)^2 \quad (\text{because of independence}) \\
&= \frac{\sigma_X^2}{N}
\end{aligned}$$

$$\text{also } E(X_i - \mu_X)(\mu_X - \hat{\mu}_X) = -E(X_i - \mu_X)^2$$

$$\therefore E \sum_{i=1}^N (X_i - \hat{\mu}_X)^2 = N\sigma^2 + \sigma^2 - 2\sigma^2 = (N-1)\sigma^2$$

$$\text{So } E\hat{\sigma}_2^2 = \frac{1}{N-1} E \sum (X_i - \hat{\mu}_X)^2 = \sigma^2$$

$$\therefore \hat{\sigma}_2^2 \text{ is an unbiased estimator of } \sigma^2.$$

Similarly sample mean is an unbiased estimator.

$$\hat{\mu}_X = \frac{1}{N} \sum_{i=1}^N X_i$$

$$E\hat{\mu}_X = \frac{1}{N} \sum_{i=1}^N E\{X_i\} = \frac{N\mu_X}{N} = \mu_X$$

#### 4.4 Variance of the estimator

The variance of the estimator  $\hat{\theta}$  is given by

$$\text{var}(\hat{\theta}) = E(\hat{\theta} - E(\hat{\theta}))^2$$

For the unbiased case

$$\text{var}(\hat{\theta}) = E(\hat{\theta} - \theta)^2$$

The variance of the estimator should be as low as possible.

An unbiased estimator  $\hat{\theta}$  is called a *minimum variance unbiased estimator* (MVUE) if

$$E(\hat{\theta} - \theta)^2 \leq E(\hat{\theta}' - \theta)^2$$

where  $\hat{\theta}'$  is any other unbiased estimator.

## 4.5 Mean square error of the estimator

$$MSE = E(\hat{\theta} - \theta)^2$$

MSE should be as small as possible. Out of all unbiased estimator, the MVUE has the minimum mean square error.

MSE is related to the bias and variance as shown below.

$$\begin{aligned} MSE &= E(\hat{\theta} - \theta)^2 = E(\hat{\theta} - E\hat{\theta} + E\hat{\theta} - \theta)^2 \\ &= E(\hat{\theta} - E\hat{\theta})^2 + E(E\hat{\theta} - \theta)^2 + 2E(\hat{\theta} - E\hat{\theta})(E\hat{\theta} - \theta) \\ &= E(\hat{\theta} - E\hat{\theta})^2 + E(E\hat{\theta} - \theta)^2 + 2(E\hat{\theta} - E\hat{\theta})(E\hat{\theta} - \theta) \\ &= \text{var}(\hat{\theta}) + b^2(\hat{\theta}) + 0 \text{ (why?)} \end{aligned}$$

So

$$MSE = \text{var}(\hat{\theta}) + b^2(\hat{\theta})$$

## 4.6 Consistent Estimators

As we have more data, the quality of estimation should be better.

This idea is used in defining the consistent estimator.

An estimator  $\hat{\theta}$  is called a consistent estimator of  $\theta$  if  $\hat{\theta}$  converges in probability to  $\theta$ .

$$\lim_{N \rightarrow \infty} P(|\hat{\theta} - \theta| \geq \varepsilon) = 0 \text{ for any } \varepsilon > 0$$

Less rigorous test is obtained by applying the Markov Inequality

$$P(|\hat{\theta} - \theta| \geq \varepsilon) \leq \frac{E(\hat{\theta} - \theta)^2}{\varepsilon^2} \quad \leftarrow \text{MSE}$$

If  $\hat{\theta}$  is an unbiased estimator ( $b(\hat{\theta}) = 0$ ), then  $MSE = \text{var}(\hat{\theta})$ .

Therefore, if

$$\lim_{N \rightarrow \infty} E(\hat{\theta} - \theta)^2 = 0, \text{ then } \hat{\theta} \text{ will be a consistent estimator.}$$

Also note that

$$MSE = \text{var}(\hat{\theta}) + b^2(\hat{\theta})$$

Therefore, if the estimator is asymptotically unbiased (i.e.  $b(\hat{\theta}) \rightarrow 0$  as  $N \rightarrow \infty$ ) and  $\text{var}(\hat{\theta}) \rightarrow 0$  as  $N \rightarrow \infty$ , then  $MSE \rightarrow 0$ .

$\therefore$  Therefore for an asymptotically unbiased estimator  $\hat{\theta}$ , if  $\text{var}(\hat{\theta}) \rightarrow 0$ , as  $N \rightarrow \infty$ , then  $\hat{\theta}$  will be a consistent estimator.



**Example 1:**  $X_1, X_2, \dots, X_N$  is an iid random sequence with unknown  $\mu_X$  and known variance  $\sigma_X^2$ .

Let  $\hat{\mu}_X = \frac{1}{N} \sum_{i=1}^N X_i$  be an estimator for  $\mu_X$ . We have already shown that  $\hat{\mu}_X$  is unbiased.

Also  $\text{var}(\hat{\mu}_X) = \frac{\sigma_X^2}{N}$  Is it a consistent estimator?

Clearly  $\lim_{N \rightarrow \infty} \text{var}(\hat{\mu}_X) = \lim_{N \rightarrow \infty} \frac{\sigma_X^2}{N} = 0$ . Therefore  $\hat{\mu}_X$  is a consistent estimator of  $\mu_X$ .

## 4.7 Sufficient Statistic

The observations  $X_1, X_2, \dots, X_N$  contain information about the unknown parameter  $\theta$ . An estimator should carry the same information about  $\theta$  as the observed data. This concept of sufficient statistic is based on this idea.

A measurable function  $\hat{\theta}(X_1, X_2, \dots, X_N)$  is called a sufficient statistic of  $\theta$  if it contains the same information about  $\theta$  as contained in the random sequence  $X_1, X_2, \dots, X_N$ . In other word the joint conditional density  $f_{X_1, X_2, \dots, X_N | \hat{\theta}(X_1, X_2, \dots, X_N)}(x_1, x_2, \dots, x_N)$  does not involve  $\theta$ .

There are a large number of sufficient statistics for a particular criterion. One has to select a sufficient statistic which has good estimation properties.

A way to check whether a statistic is sufficient or not is through the *Factorization theorem* which states:

$\hat{\theta}(X_1, X_2, \dots, X_N)$  is a sufficient statistic of  $\theta$  if

$$f_{X_1, X_2, \dots, X_N | \theta}(x_1, x_2, \dots, x_N) = g(\theta, \hat{\theta})h(x_1, x_2, \dots, x_N)$$

where  $g(\theta, \hat{\theta})$  is a non-constant and nonnegative function of  $\theta$  and  $\hat{\theta}$  and  $h(x_1, x_2, \dots, x_N)$  does not involve  $\theta$  and is a nonnegative function of  $x_1, x_2, \dots, x_N$ .

**Example 2:** Suppose  $X_1, X_2, \dots, X_N$  is an iid Gaussian sequence with unknown mean  $\mu_X$  and known variance 1.

Then  $\hat{\mu}_X = \frac{1}{N} \sum_{i=1}^N X_i$  is a sufficient statistic of.

$$\begin{aligned}
f_{X_1, X_2, \dots, X_N / \mu_X}(x_1, x_2, \dots, x_N) &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i - \mu_X)^2} \\
&= \frac{1}{(\sqrt{2\pi})^N} e^{-\frac{1}{2} \sum_{i=1}^N (x_i - \mu_X)^2} \\
&= \frac{1}{(\sqrt{2\pi})^N} e^{-\frac{1}{2} \sum_{i=1}^N (x_i - \hat{\mu} + \hat{\mu} - \mu)^2} \\
&= \frac{1}{(\sqrt{2\pi})^N} e^{-\frac{1}{2} \sum_{i=1}^N ((x_i - \hat{\mu})^2 + (\hat{\mu} - \mu_X)^2 + 2(x_i - \hat{\mu})(\hat{\mu} - \mu_X))} \\
&= \frac{1}{(\sqrt{2\pi})^N} e^{-\frac{1}{2} \sum_{i=1}^N (x_i - \hat{\mu})^2} e^{-\frac{1}{2} \sum_{i=1}^N ((\hat{\mu} - \mu_X)^2)} e^0 \text{ (why?)}
\end{aligned}$$

Because

The first exponential is a function of  $x_1, x_2, \dots, x_N$  and the second exponential is a function of  $\mu_X$  and  $\hat{\mu}_X$ . Therefore  $\hat{\mu}_X$  is a sufficient statistics of  $\mu_X$ .

#### 4.8 Cramer Rao theorem

We described about the *Minimum Variance Unbiased Estimator (MVUE)* which is a very good estimator

$\hat{\theta}$  is an MVUE if

$$E(\hat{\theta}) = \theta$$

and  $Var(\hat{\theta}) \leq Var(\hat{\theta}')$

where  $\hat{\theta}'$  is any other unbiased estimator of  $\theta$ .

Can we reduce the variance of an unbiased estimator indefinitely? The answer is given by the Cramer Rao theorem.

Suppose  $\hat{\theta}$  is an unbiased estimator of random sequence. Let us denote the sequence by the vector

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$$

Let  $f_{\mathbf{X}}(x_1, \dots, x_N / \theta)$  be the joint density function which characterises  $\mathbf{X}$ . This function is also called likelihood function.  $\theta$  may also be random. In that case likelihood function will represent conditional joint density function.

$L(\mathbf{x} / \theta) = \ln f_{\mathbf{X}}(x_1, \dots, x_N / \theta)$  is called log likelihood function.

## 4.9 Statement of the Cramer Rao theorem

If  $\hat{\theta}$  is an unbiased estimator of  $\theta$ , then

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)}$$

where  $I(\theta) = E\left(\frac{\partial L}{\partial \theta}\right)^2$  and  $I(\theta)$  is a measure of average information in the random sequence and is called Fisher information statistic.

The equality of CR bound holds if  $\frac{\partial L}{\partial \theta} = c(\hat{\theta} - \theta)$  where  $c$  is a constant.

**Proof:**  $\hat{\theta}$  is an unbiased estimator of  $\theta$

$$\therefore E(\hat{\theta} - \theta) = 0.$$

$$\Rightarrow \int_{-\infty}^{\infty} (\hat{\theta} - \theta) f_{\mathbf{x}}(\mathbf{x} / \theta) d\mathbf{x} = 0.$$

Differentiate with respect to  $\theta$ , we get

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} \{(\hat{\theta} - \theta) f_{\mathbf{x}}(\mathbf{x} / \theta)\} d\mathbf{x} = 0.$$

(Since line of integration are not function of  $\theta$ .)

$$= \int_{-\infty}^{\infty} (\hat{\theta} - \theta) \frac{\partial}{\partial \theta} f_{\mathbf{x}}(\mathbf{x} / \theta) d\mathbf{x} - \int_{-\infty}^{\infty} f_{\mathbf{x}}(\mathbf{x} / \theta) d\mathbf{x} = 0.$$

$$\therefore \int_{-\infty}^{\infty} (\hat{\theta} - \theta) \frac{\partial}{\partial \theta} f_{\mathbf{x}}(\mathbf{x} / \theta) d\mathbf{x} = \int_{-\infty}^{\infty} f_{\mathbf{x}}(\mathbf{x} / \theta) d\mathbf{x} = 1. \quad (1)$$

Note that  $\frac{\partial}{\partial \theta} f_{\mathbf{x}}(\mathbf{x} / \theta) = \frac{\partial}{\partial \theta} \{\ln f_{\mathbf{x}}(\mathbf{x} / \theta)\} f_{\mathbf{x}}(\mathbf{x} / \theta)$

$$= \left(\frac{\partial L}{\partial \theta}\right) f_{\mathbf{x}}(\mathbf{x} / \theta)$$

Therefore, from (1)

$$\int_{-\infty}^{\infty} (\hat{\theta} - \theta) \left\{ \frac{\partial}{\partial \theta} L(\mathbf{x} / \theta) \right\} f_{\mathbf{x}}(\mathbf{x} / \theta) d\mathbf{x} = 1.$$

So that

$$\left\{ \int_{-\infty}^{\infty} (\hat{\theta} - \theta) \sqrt{f_{\mathbf{x}}(\mathbf{x} / \theta)} \frac{\partial}{\partial \theta} L(\mathbf{x} / \theta) \sqrt{f_{\mathbf{x}}(\mathbf{x} / \theta)} d\mathbf{x} \right\}^2 = 1. \quad (2)$$

since  $f_{\mathbf{x}}(\mathbf{x} / \theta)$  is  $\geq 0$ .

Recall the Cauchy Schwarz Inequality

$$|\langle \mathbf{a}, \mathbf{b} \rangle|^2 \leq \|\mathbf{a}\|^2 \|\mathbf{b}\|^2$$

where the equality holds when  $\mathbf{a} = c\mathbf{b}$  ( where c is any scalar ).

Applying this inequality to the L.H.S. of equation (2) we get

$$\begin{aligned} & \left( \int_{-\infty}^{\infty} (\hat{\theta} - \theta) \sqrt{f_{\mathbf{x}}(\mathbf{x}/\theta)} \frac{\partial}{\partial \theta} L(\mathbf{x}/\theta) \sqrt{f_{\mathbf{x}}(\mathbf{x}/\theta)} d\mathbf{x} \right)^2 \\ & \leq \int_{-\infty}^{\infty} (\hat{\theta} - \theta)^2 f_{\mathbf{x}}(\mathbf{x}/\theta) d\mathbf{x} \int_{-\infty}^{\infty} \left( \frac{\partial}{\partial \theta} L(\mathbf{x}/\theta) \right)^2 f_{\mathbf{x}}(\mathbf{x}/\theta) d\mathbf{x} \\ & = \text{var}(\hat{\theta}) I(\theta) \\ & \therefore L.H.S \leq \text{var}(\hat{\theta}) I(\theta) \end{aligned}$$

But R.H.S. = 1

$$\text{var}(\hat{\theta}) I(\theta) \geq 1.$$

$$\therefore \text{var}(\hat{\theta}) \geq \frac{1}{I(\theta)},$$

which is the Cramer Rao Inequality.

The equality will hold when

$$\frac{\partial}{\partial \theta} \{L(\mathbf{x}/\theta) \sqrt{f_{\mathbf{x}}(\mathbf{x}/\theta)}\} = c(\hat{\theta} - \theta) \sqrt{f_{\mathbf{x}}(\mathbf{x}/\theta)},$$

so that

$$\boxed{\frac{\partial L(\mathbf{x}/\theta)}{\partial \theta} = c(\hat{\theta} - \theta)}$$

Also from  $\int_{-\infty}^{\infty} f_{\mathbf{x}}(\mathbf{x}/\theta) d\mathbf{x} = 1$ , we get

$$\begin{aligned} & \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_{\mathbf{x}}(\mathbf{x}/\theta) d\mathbf{x} = 0 \\ & \therefore \int_{-\infty}^{\infty} \frac{\partial L}{\partial \theta} f_{\mathbf{x}}(\mathbf{x}/\theta) d\mathbf{x} = 0 \end{aligned}$$

Taking the partial derivative with respect to  $\theta$  again, we get

$$\begin{aligned} & \int_{-\infty}^{\infty} \left\{ \frac{\partial^2 L}{\partial \theta^2} f_{\mathbf{x}}(\mathbf{x}/\theta) + \frac{\partial L}{\partial \theta} \frac{\partial}{\partial \theta} f_{\mathbf{x}}(\mathbf{x}/\theta) \right\} d\mathbf{x} = 0 \\ & \therefore \int_{-\infty}^{\infty} \left\{ \frac{\partial^2 L}{\partial \theta^2} f_{\mathbf{x}}(\mathbf{x}/\theta) + \left( \frac{\partial L}{\partial \theta} \right)^2 f_{\mathbf{x}}(\mathbf{x}/\theta) \right\} d\mathbf{x} = 0 \\ & E \left( \frac{\partial L}{\partial \theta} \right)^2 = - E \frac{\partial^2 L}{\partial \theta^2} \end{aligned}$$

If  $\hat{\theta}$  satisfies CR -bound with equality, then  $\hat{\theta}$  is called an efficient estimator.

**Remark:**

- (1) If the information  $I(\theta)$  is more, the variance of the estimator  $\hat{\theta}$  will be less.
- (2) Suppose  $X_1, \dots, X_N$  are iid. Then

$$I_1(\theta) = E \left( \frac{\partial}{\partial \theta} \ln(f_{X_1/\theta}(x)) \right)^2$$

$$\therefore I_N(\theta) = E \left( \frac{\partial}{\partial \theta} \ln(f_{X_1, X_2, \dots, X_N/\theta}(x_1, x_2, \dots, x_N)) \right)^2$$

$$= NI_1(\theta)$$

**Example 3:**

Let  $X_1, \dots, X_N$  are iid Gaussian random sequence with known variance  $\sigma^2$  and unknown mean  $\mu$ .

Suppose  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$  which is unbiased.

Find CR bound and hence show that  $\hat{\mu}$  is an efficient estimator.

Likelihood function

$f_{\mathbf{X}}(x_1, x_2, \dots, x_N / \theta)$  will be product of individual densities (since iid)

$$\therefore f_{\mathbf{X}}(x_1, x_2, \dots, x_N / \theta) = \frac{1}{(\sqrt{2\pi})^N \sigma^N} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2}$$

$$\text{so that } L(\mathbf{X} / \mu) = -\ln(\sqrt{2\pi})^N \sigma^N - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\text{Now } \frac{\partial L}{\partial \mu} = 0 - \frac{1}{2\sigma^2} (-2) \sum_{i=1}^N (X_i - \mu)$$

$$\therefore \frac{\partial^2 L}{\partial \mu^2} = -\frac{N}{\sigma^2}$$

$$\text{So that } E \frac{\partial^2 L}{\partial \mu^2} = -\frac{N}{\sigma^2}$$

$$\therefore \text{CR Bound} = \frac{1}{I(\theta)} = \frac{1}{-E \frac{\partial^2 L}{\partial \mu^2}} = \frac{1}{\frac{N}{\sigma^2}} = \frac{\sigma^2}{N}$$

$$\frac{\partial L}{\partial \theta} = \frac{1}{2\sigma^2} \sum_{i=1}^N (X_i - \mu) = \frac{N}{\sigma^2} \left( \sum_i \frac{X_i}{N} - \mu \right)$$

$$= \frac{N}{\sigma^2} (\hat{\mu} - \mu)$$

$$\text{Hence } -\frac{\partial L}{\partial \theta} = c(\hat{\theta} - \theta)$$

and  $\hat{\mu}$  is an efficient estimator.

#### 4.10 Criteria for Estimation

The estimation of a parameter is based on several well-known criteria. Each of the criteria tries to optimize some functions of the observed samples with respect to the unknown parameter to be estimated. Some of the most popular estimation criteria are:

- Maximum Likelihood
- Minimum Mean Square Error.
- Baye's Method.
- Maximum Entropy Method.

#### 4.11 Maximum Likelihood Estimator (MLE)

Given a random sequence  $X_1, \dots, X_N$  and the joint density function  $f_{X_1, \dots, X_N / \theta}(x_1, x_2, \dots, x_N)$  which depends on an unknown nonrandom parameter  $\theta$ .

$f_{\mathbf{x}}(x_1, x_2, \dots, x_N / \theta)$  is called the likelihood function (for continuous function ..., for discrete it will be joint probability mass function).

$L(\mathbf{x} / \theta) = \ln f_{\mathbf{x}}(x_1, x_2, \dots, x_N / \theta)$  is called log likelihood function.

The maximum likelihood estimator  $\hat{\theta}_{MLE}$  is such an estimator that

$$f_{\mathbf{x}}(x_1, x_2, \dots, x_N / \hat{\theta}_{MLE}) \geq f_{\mathbf{x}}(x_1, x_2, \dots, x_N / \theta), \forall \theta$$

If the likelihood function is differentiable w.r.t.  $\theta$ , then  $\hat{\theta}_{MLE}$  is given by

$$\frac{\partial}{\partial \theta} f_{\mathbf{x}}(x_1, \dots, x_N / \theta) \Big|_{\hat{\theta}_{MLE}} = 0$$

$$\text{or } \frac{\partial L(\mathbf{x} | \theta)}{\partial \theta} \Big|_{\hat{\theta}_{MLE}} = 0$$

Thus the MLE is given by the solution of the likelihood equation given above.

If we have a number of unknown parameters given by  $\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_N \end{bmatrix}$

Then MLE is given by a set of conditions.

$$\left. \frac{\partial L}{\partial \theta_1} \right|_{\theta_1 = \hat{\theta}_{1MLE}} = \left. \frac{\partial L}{\partial \theta_2} \right|_{\theta_2 = \hat{\theta}_{2MLE}} = \dots = \left. \frac{\partial L}{\partial \theta_M} \right|_{\theta_M = \hat{\theta}_{MMLE}} = 0$$

#### **Example 4:**

Let  $X_1, \dots, X_N$  are independent identically distributed sequence of  $N(\mu, \sigma^2)$  distributed random variables. Find MLE for  $\mu, \sigma^2$ .

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_N / \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x_i - \mu}{\sigma}\right)^2}$$

$$L(X / \mu, \sigma^2) = \ln f_{\mathbf{X}}(X_1, \dots, X_N / \mu, \sigma^2)$$

$$= -N \ln \frac{1}{\sqrt{2\pi}\sigma} - N \ln \sigma - \frac{1}{2} \sum_{i=1}^N \left( \frac{x_i - \mu}{\sigma} \right)^2$$

$$\frac{\partial L}{\partial \mu} = 0 \Rightarrow \sum_{i=1}^N \left( \frac{x_i - \mu}{\sigma} \right)^2 = \sum_{i=1}^N (x_i - \hat{\mu}_{MLE}) = 0$$

$$\frac{\partial L}{\partial \sigma} = 0 = -\frac{N}{\hat{\sigma}_{MLE}} + \frac{\sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2}{\hat{\sigma}_{MLE}^3} = 0$$

Solving we get

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i \quad \text{and}$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$$

#### **Example 5:**

Let  $X_1, \dots, X_N$  are independent identically distributed sequence with

$$f_{\mathbf{X}/\theta}(x) = \frac{1}{2} e^{-|x-\theta|} \quad -\infty < x < \infty$$

Show that the median of  $X_1, \dots, X_N$  is the MLE for  $\theta$ .

$$f_{X_1, X_2, \dots, X_N / \theta}(x_1, x_2, \dots, x_N) = \frac{1}{2^N} e^{-\sum_{i=1}^N |x_i - \theta|}$$

$$L(X / \theta) = \ln f_{\mathbf{X}/\theta}(x_1, \dots, x_N) = -N \ln 2 - \sum_{i=1}^N |x_i - \theta|$$

$$\sum_{i=1}^N |x_i - \theta| \quad \text{is minimized by} \quad \text{median}(X_1, \dots, X_N)$$

### Some properties of MLE (without proof)

- MLE may be biased or unbiased, asymptotically unbiased.
- MLE is consistent estimator.
- If an efficient estimator exists, it is the MLE estimator.

An efficient estimator  $\hat{\theta}$  exists  $\Rightarrow$

$$\frac{\partial}{\partial \theta} L(\mathbf{x} / \theta) = c(\hat{\theta} - \theta)$$

at  $\theta = \hat{\theta}$ ,

$$\left. \frac{\partial L(\mathbf{x} / \theta)}{\partial \theta} \right|_{\hat{\theta}} = c(\hat{\theta} - \hat{\theta}) = 0$$

$\Rightarrow \hat{\theta}$  is the MLE estimator.

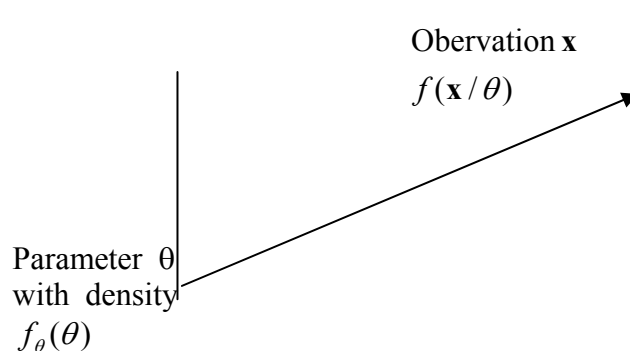
- Invariance Properties of MLE:

If  $\hat{\theta}_{MLE}$  is the MLE of  $\theta$ , then  $h(\hat{\theta}_{MLE})$  is the MLE of  $h(\theta)$ , where  $h(\theta)$  is an invertible function of  $\theta$ .

## 4.12 Bayescan Estimators

We may have some prior information about  $\theta$  in a sense that some values of  $\theta$  are more likely (*a priori* information). We can represent this prior information in the form of a prior density function.

In the following we omit the suffix in density functions just for notational simplicity.



The likelihood function will now be the conditional density  $f(\mathbf{x} / \theta)$ .

$$f_{\mathbf{x}, \theta}(\mathbf{x}, \theta) = f_{\theta}(\theta) f_{\mathbf{x} / \theta}(\mathbf{x})$$

Also we have the Bayes rule

$$f_{\theta / \mathbf{x}}(\theta) = \frac{f_{\theta}(\theta) f_{\mathbf{x} / \theta}(\mathbf{x})}{f_{\mathbf{x}}(\mathbf{x})}$$

where  $f_{\theta / \mathbf{x}}(\theta)$  is the *a posteriori* density function



The parameter  $\theta$  is a random variable and the estimator  $\hat{\theta}(\mathbf{x})$  is another random variable.

Estimation error  $\varepsilon = \hat{\theta} - \theta$ .

We associate a cost function  $C(\hat{\theta}, \theta)$  with every estimator  $\hat{\theta}$ . It represents the positive penalty with each wrong estimation.

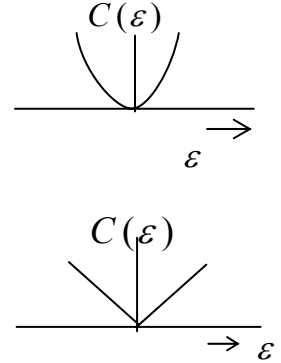
Thus  $C(\hat{\theta}, \theta)$  is a non negative function.

The three most popular cost functions are:

**Quadratic cost function**  $(\hat{\theta} - \theta)^2$

**Absolute cost function**  $|\hat{\theta} - \theta|$

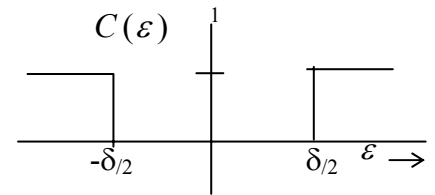
**Hit or miss cost function** (also called uniform cost function)  
minimising means minimising on an average)



#### 4.13 Bayesean Risk function or average cost

$$\bar{C} = EC(\theta, \hat{\theta}) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} C(\theta, \hat{\theta}) f_{\mathbf{x}, \theta}(\mathbf{x}, \theta) d\mathbf{x} d\theta$$

The estimator seeks to minimize the Bayesean Risk.



##### Case I. Quadratic Cost Function

$$C = (\theta, \hat{\theta}) = (\hat{\theta} - \theta)^2$$

Estimation problem is

$$\text{Minimize} \quad \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\theta - \hat{\theta})^2 f_{\mathbf{x}, \theta}(\mathbf{x}, \theta) d\mathbf{x} d\theta$$

with respect to  $\hat{\theta}$ .

This is equivalent to minimizing

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\theta - \hat{\theta})^2 f(\theta | \mathbf{x}) f(\mathbf{x}) d\theta d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} (\theta - \hat{\theta})^2 f(\theta | \mathbf{x}) d\theta \right) f(\mathbf{x}) d\mathbf{x} \end{aligned}$$

Since  $f(\mathbf{x})$  is always +ve, the above integral will be minimum if the inner integral is minimum. This results in the problem:

$$\text{Minimize} \quad \int_{-\infty}^{\infty} (\theta - \hat{\theta})^2 f(\theta | \mathbf{x}) d\theta$$

with respect to  $\hat{\theta}$ .

$$\Rightarrow \frac{\partial}{\partial \hat{\theta}} \int_{-\infty}^{\infty} (\hat{\theta} - \theta)^2 f_{\Theta/\mathbf{X}}(\theta) d\theta = 0$$

$$\Rightarrow -2 \int_{-\infty}^{\infty} (\hat{\theta} - \theta) f_{\Theta/\mathbf{X}}(\theta) d\theta = 0$$

$$\Rightarrow \hat{\theta} \int_{-\infty}^{\infty} f_{\Theta/\mathbf{X}}(\theta) d\theta = \int_{-\infty}^{\infty} \theta f_{\Theta/\mathbf{X}}(\theta) d\theta$$

$$\Rightarrow \hat{\theta} = \int_{-\infty}^{\infty} \theta f_{\Theta/\mathbf{X}}(\theta) d\theta$$

$\therefore \hat{\theta}$  is the conditional mean or mean of the a posteriori density. Since we are minimizing quadratic cost it is also called *minimum mean square error estimator* (MMSE).

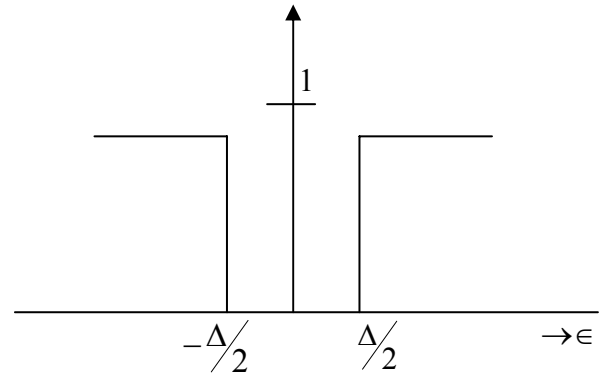
### Salient Points

- Information about distribution of  $\theta$  available.
- *a priori* density function  $f_{\Theta}(\theta)$  is available. This denotes how observed data depend on  $\theta$
- We have to determine a posteriori density  $f_{\Theta/\mathbf{X}}(\theta)$ . This is determined from the Bayes rule.

### Case II Hit or Miss Cost Function

$$\text{Risk } \bar{C} = EC(\theta, \hat{\theta}) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} C(\theta, \hat{\theta}) f_{\mathbf{X}, \Theta}(\mathbf{x}, \theta) d\mathbf{x} d\theta$$

$$\begin{aligned} &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} c(\theta, \hat{\theta}) f_{\Theta/\mathbf{X}}(\theta) f_{\mathbf{X}}(\mathbf{x}) d\theta d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} c(\theta, \hat{\theta}) f_{\Theta/\mathbf{X}}(\theta) d\theta \right) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$



We have to minimize

$$\int_{-\infty}^{\infty} C(\theta, \hat{\theta}) f_{\Theta/\mathbf{X}}(\theta) d\theta \quad \text{with respect to } \hat{\theta}.$$

This is equivalent to minimizing

$$= 1 - \int_{\hat{\theta} - \frac{\Delta}{2}}^{\hat{\theta} + \frac{\Delta}{2}} f_{\Theta/\mathbf{X}}(\theta) d\theta$$

This minimization is equivalent to maximization of

$$\int_{\hat{\theta} - \frac{\Delta}{2}}^{\hat{\theta} + \frac{\Delta}{2}} f_{\Theta/X}(\theta) d\theta \cong \Delta f_{\Theta/X}(\hat{\theta}) \quad \text{when } \Delta \text{ is very small}$$

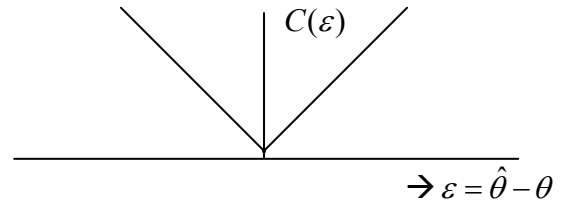
This will be maximum if  $f_{\Theta/X}(\theta)$  is maximum. That means select that value of  $\hat{\theta}$  that maximizes the a posteriori density. So this is known as maximum *a posteriori* estimation (MAP) principle.

This estimator is denoted by  $\hat{\theta}_{MAP}$ .

### Case III

$$C(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$$

$$\begin{aligned} \bar{C} &= \text{Average cost} = E|\hat{\theta} - \theta| \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |\hat{\theta} - \theta| f_{\theta, \mathbf{x}}(\theta, \mathbf{x}) d\theta d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |\hat{\theta} - \theta| f_{\theta}(\theta) f_{\mathbf{x}/\theta}(\mathbf{x}) d\theta d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |\hat{\theta} - \theta| f_{\mathbf{x}/\theta}(\mathbf{x}) d\mathbf{x} f_{\theta}(\theta) d\theta \end{aligned}$$



For the minimum

$$\begin{aligned} \frac{\partial}{\partial \hat{\theta}} \int_{-\infty}^{\infty} C(\hat{\theta}, \theta) f_{\theta/X}(\theta | x) d\theta &= 0 \\ \frac{\partial}{\partial \hat{\theta}} \left\{ \int_{-\infty}^{\hat{\theta}} (\hat{\theta} - \theta) f_{\theta/X}(\theta | x) d\theta + \int_{\hat{\theta}}^{\infty} (\theta - \hat{\theta}) f_{\theta/X}(\theta | x) d\theta \right\} &= 0 \end{aligned}$$

Leibniz rule for differentiation of integration

$$\begin{aligned} \frac{\partial}{\partial u} \int_{\theta_1(u)}^{\theta_2(u)} h(u, v) dv &= \int_{\theta_1(u)}^{\theta_2(u)} \frac{\partial h(u, v)}{\partial u} dv \\ &+ \frac{d\theta_2(u)}{du} h(u, \theta_2(u)) - \frac{d\theta_1(u)}{du} h(u, \theta_1(u)) \end{aligned}$$

Applying Leibniz rule we get

$$\int_{-\infty}^{\hat{\theta}} f_{\theta/X}(\theta | x) d\theta - \int_{\hat{\theta}}^{\infty} f_{\theta/X}(\theta | x) d\theta = 0$$

At the  $\hat{\theta}_{MAE}$

$$\int_{-\infty}^{\hat{\theta}_{MAE}} f_{\theta/X}(\theta | x) d\theta - \int_{\hat{\theta}_{MAE}}^{\infty} f_{\theta/X}(\theta | x) d\theta = 0$$

So  $\hat{\theta}_{MAE}$  is the median of the *a posteriori* density

**Example 6:**

Let  $X_1, X_2, \dots, X_N$  be an *iid* Gaussian sequence with unity Variance and unknown mean  $\theta$ . Further  $\theta$  is known to be a 0-mean Gaussian with Unity Variance. Find the MAP estimator for  $\theta$ .

**Solution:** We are given

$$f_{\theta}(\theta) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\theta^2}$$

$$f_{\theta/\mathbf{x}}(\theta) = \frac{1}{(\sqrt{2\pi})^N} e^{-\frac{\sum_{i=1}^N (x_i - \theta)^2}{2}}$$

$$\text{Therefore } f_{\theta/\mathbf{x}}(\theta) = \frac{f_{\theta}(\theta)f_{\mathbf{x}/\theta}(\mathbf{x})}{f_{\mathbf{x}}(\mathbf{x})}$$

We have to find  $\theta$ , such that  $f_{\theta/\mathbf{x}}(\theta)$  is maximum.

Now  $f_{\theta/\mathbf{x}}(\theta)$  is maximum when  $f_{\theta}(\theta)f_{\mathbf{x}/\theta}(\mathbf{x})$  is maximum.

$\Rightarrow \ln f_{\theta}(\theta)f_{\mathbf{x}/\theta}(\mathbf{x})$  is maximum

$$\Rightarrow -\frac{1}{2}\theta^2 - \sum_{i=1}^N \frac{(x_i - \theta)^2}{2} \text{ is maximum}$$

$$\Rightarrow \theta - \sum_{i=1}^N (x_i - \theta) \Big|_{\theta=\hat{\theta}_{MAP}} = 0$$

$$\Rightarrow \hat{\theta}_{MAP} = \frac{1}{N+1} \sum_{i=1}^N x_i$$

**Example 7:**

Consider single observation  $X$  that depends on a random parameter  $\theta$ . Suppose  $\theta$  has a prior distribution

$$f_{\theta}(\theta) = \lambda e^{-\lambda\theta} \quad \text{for } \theta \geq 0, \quad \lambda > 0$$

and

$$f_{X/\theta}(x) = \theta e^{-\theta x} \quad |x| > 0$$

find the MAP estimation for  $\theta$ .

$$f_{\theta/X}(\theta) = \frac{f_{\theta}(\theta)f_{X/\theta}(x)}{f_X(x)}$$

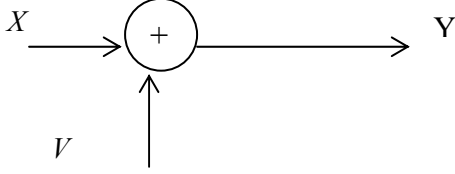
$$\ln(f(\theta | x)) = \ln(f_{\theta}(\theta)) + \ln(f_{X/\theta}(x)) - \ln f_X(x)$$

Therefore MAP estimator is given by.

$$\frac{\partial}{\partial \theta} \ln f_{X/\Theta}(x) \Big|_{\hat{\theta}_{\text{MAP}}} = 0$$

$$\Rightarrow \hat{\theta}_{\text{MAP}} = \frac{1}{\lambda + X}$$

**Example 8:** Binary Communication problem



$X$  is a binary random variable with

$$X = \begin{cases} 1 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/2 \end{cases}$$

$V$  is the Gaussian noise with mean 0 and variance  $\sigma^2$ .

To find the MSE for  $X$  from the observed data  $Y$ .

Then

$$f_X(x) = \frac{1}{2} [\delta(x-1) + \delta(x+1)]$$

$$f_{Y/X}(y/x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(y-x)^2/2\sigma^2}$$

$$f_{X/Y}(x/y) = \frac{f_X(x)f_{Y/X}(y/x)}{\int_{-\infty}^{\infty} f_X(x)f_{Y/X}(y/x)dx}$$

$$= \frac{e^{-(y-x)^2/2\sigma^2} [\delta(x-1) + \delta(x+1)]}{e^{-(y-1)^2/2\sigma^2} + e^{-(y+1)^2/2\sigma^2}}$$

$$\hat{X}_{MMSE} = E(X/Y) = \int_{-\infty}^{\infty} x \frac{e^{-(y-x)^2/2\sigma^2} [\delta(x-1) + \delta(x+1)]}{e^{-(y-1)^2/2\sigma^2} + e^{-(y+1)^2/2\sigma^2}} dx$$

Hence

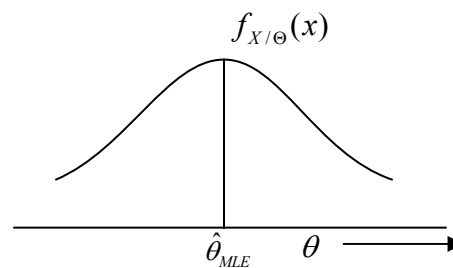
$$= \frac{e^{-(y-1)^2/2\sigma^2} - e^{-(y+1)^2/2\sigma^2}}{e^{-(y-1)^2/2\sigma^2} + e^{-(y+1)^2/2\sigma^2}}$$

$$= \tanh(y/\sigma^2)$$

## To summarise

**MLE:**

Simplest



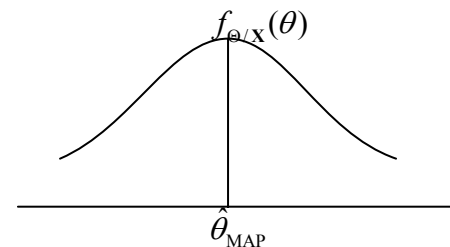
**MMSE:**

$$\hat{\theta}_{MMSE} = E(\Theta / \mathbf{X})_{MMSE}$$

- Find a posteriori density.
- Find the average value by integration
- Lots of calculation hence it is computationally exhaustive.

**MAP:**

$\hat{\theta}_{MAP}$  = Mode of the *a posteriori* density  $f_{\Theta/\mathbf{X}}(\theta)$ .



### 4.14 Relation between $\hat{\theta}_{MAP}$ and $\hat{\theta}_{MLE}$

From

$$f_{\Theta/\mathbf{X}}(\theta) = \frac{f_{\Theta}(\theta)f_{\mathbf{X}/\Theta}(\mathbf{x})}{f_{\mathbf{X}}(\mathbf{x})}$$

$$\ln(f_{\Theta/\mathbf{X}}(\theta)) = \ln(f_{\Theta}(\theta)) + \ln(f_{\mathbf{X}/\Theta}(\mathbf{x})) - \ln(f_{\mathbf{X}}(\mathbf{x}))$$

$\hat{\theta}_{MAP}$  is given by

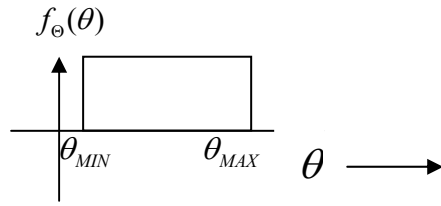
$$\frac{\partial}{\partial \theta} \ln f_{\Theta}(\theta) + \underbrace{\frac{\partial}{\partial \theta} \ln(f_{\mathbf{X}/\Theta}(\mathbf{x}))}_{\text{likelihood function}} = 0$$

a priori density      likelihood function.

Suppose  $\theta$  is uniformly distributed between  $\theta_{MIN}$  and  $\theta_{MAX}$ .

Then

$$\frac{\partial}{\partial \theta} \ln f_{\Theta}(\theta) = 0$$



If  $\theta_{MIN} \leq \hat{\theta}_{MLE} \leq \theta_{MAX}$   
then  $\hat{\theta}_{MAP} = \hat{\theta}_{MLE}$

If  $\hat{\theta}_{MAP} \leq \theta_{MIN}$   
then  $\hat{\theta}_{MAP} = \theta_{MIN}$

If  $\hat{\theta}_{MLE} \geq \theta_{MAX}$   
then  $\hat{\theta}_{MAP} = \theta_{MAX}$

## **SECTION – III**

### **OPTIMAL FILTERING**

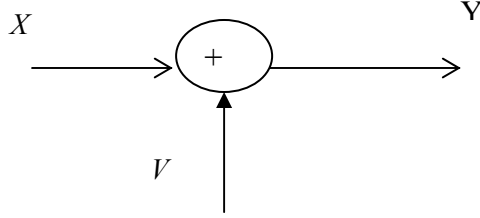


## CHAPTER – 5: WIENER FILTER

### 5.1 Estimation of signal in presence of white Gaussian noise (WGN)

Consider the signal model is

$$Y[n] = X[n] + V[n]$$



where  $Y[n]$  is the observed signal,  $X[n]$  is 0-mean Gaussian with variance 1 and  $V[n]$  is a white Gaussian sequence mean 0 and variance 1. The problem is to find the best guess for  $X[n]$  given the observation  $Y[i], i = 1, 2, \dots, n$

Maximum likelihood estimation for  $X[n]$  determines that value of  $X[n]$  for which the sequence  $Y[i], i = 1, 2, \dots, n$  is most likely. Let us represent the random sequence

$$\mathbf{Y}[\mathbf{n}] = [Y[n], Y[n-1], \dots, Y[1]]'$$

$Y[i], i = 1, 2, \dots, n$  by the random vector and the value sequence  $y[1], y[2], \dots, y[n]$  by

$$\mathbf{y}[\mathbf{n}] = [y[n], y[n-1], \dots, y[1]]'$$

The likelihood function  $f_{\mathbf{Y}[\mathbf{n}]/X[n]}(\mathbf{y}[\mathbf{n}]/x[n])$  will be Gaussian with mean  $x[n]$

$$f_{\mathbf{Y}[\mathbf{n}]/X[n]}(\mathbf{y}[\mathbf{n}]/x[n]) = \frac{1}{(\sqrt{2\pi})^n} e^{-\sum_{i=1}^n \frac{(y[i]-x[n])^2}{2}}$$

Maximum likelihood will be given by

$$\frac{\partial}{\partial x[n]} (f_{\mathbf{Y}[\mathbf{n}]/X[n]}(\mathbf{y}[\mathbf{n}]/x[n])) \Big|_{\hat{x}_{MLE}[n]} = 0$$

$$\Rightarrow \hat{x}_{MLE}[n] = \frac{1}{n} \sum_{i=1}^n y[i]$$

Similarly, to find  $\hat{x}_{MAP}[n]$  and  $\hat{x}_{MMSE}[n]$  we have to find *a posteriori density*

$$\begin{aligned} f_{X[n]/\mathbf{Y}[\mathbf{n}]}(x[n]/\mathbf{y}[\mathbf{n}]) &= \frac{f_{X[n]}(x[n]) f_{\mathbf{Y}[\mathbf{n}]/X[n]}(\mathbf{y}[\mathbf{n}]/x[n])}{f_{\mathbf{Y}[\mathbf{n}]}(\mathbf{y}[\mathbf{n}])} \\ &= \frac{1}{f_{\mathbf{Y}[\mathbf{n}]}(\mathbf{y}[\mathbf{n}])} e^{-\frac{1}{2}x^2[n] - \sum_{i=1}^n \frac{(y[i]-x[n])^2}{2}} \end{aligned}$$

Taking logarithm

$$\log_e f_{X[n]/Y[n]}(x[n]) = -\frac{1}{2}x^2[n] - \sum_{i=1}^n \frac{(y[i] - x[n])^2}{2} - \log_e f_{Y[n]}(y[n])$$

$\log_e f_{X[n]/Y[n]}(x[n])$  is maximum at  $\hat{x}_{MAP}[n]$ . Therefore, taking partial derivative of  $\log_e f_{X[n]/Y[n]}(x[n])$  with respect to  $x[n]$  and equating it to 0, we get

$$x[n] - \sum_{i=1}^n (y[i] - x[n]) \Big|_{\hat{x}_{MAP}[n]} = 0$$

$$\hat{x}_{MAP}[n] = \frac{\sum_{i=1}^n y[i]}{n+1}$$

Similarly the minimum mean-square error estimator is given by

$$\hat{x}_{MMSE}[n] = E(X[n]/y[n]) = \frac{\sum_{i=1}^n y[i]}{n+1}$$

- For MMSE we have to know the joint probability structure of the channel and the source and hence the *a posteriori pdf*.
- Finding pdf is computationally very exhaustive and nonlinear.
- Normally we may be having the estimated values first-order and second-order statistics of the data

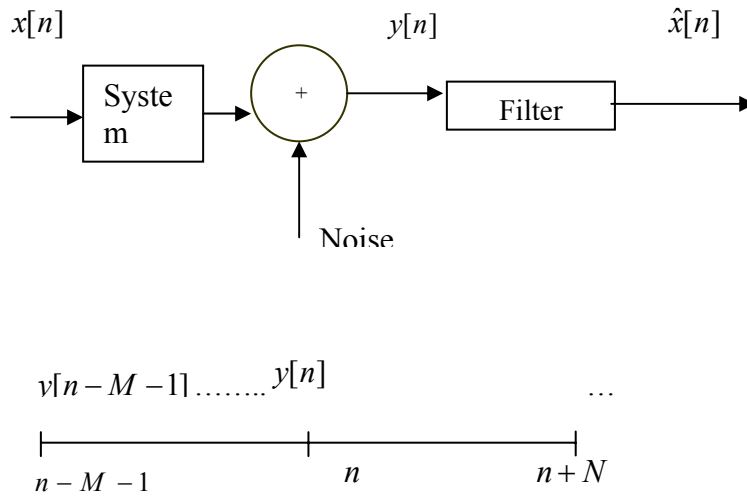
We look for a simpler estimator.

The answer is ***Optimal filtering or Wiener filtering***

We have seen that we can estimate an unknown signal (*desired signal*)  $x[n]$  from an observed signal  $y[n]$  on the basis of the known joint distributions of  $y[n]$  and  $x[n]$ . We could have used the criteria like MMSE or MAP that we have applied for parameter estimations. But such estimations are generally non-linear, require the computation of *a posteriori* probabilities and involves computational complexities.

The approach taken by Wiener is to specify a form for the estimator that depends on a number of parameters. The minimization of errors then results in determination of an optimal set of estimator parameters. A mathematically simple and computationally easier estimator is obtained by assuming a linear structure for the estimator.

## 5.2 Linear Minimum Mean Square Error Estimator



The linear minimum mean square error criterion is illustrated in the above figure. The problem can be stated as follows:

Given observations of data  $y[n-M+1], y[n-M+2], \dots, y[n], \dots, y[n+N]$ , determine a set of parameters  $h[M-1], h[M-2], \dots, h[0], \dots, h[-N]$  such that

$$\hat{x}[n] = \sum_{i=-N}^{M-1} h[i]y[n-i]$$

and the mean square error  $E(x[n] - \hat{x}[n])^2$  is a minimum with respect to  $h[-N], h[-N+1], \dots, h[0], h[1], \dots, h[M-1]$ .

This minimization problem results in an elegant solution if we assume joint stationarity of the signals  $x[n]$  and  $y[n]$ . The estimator parameters can be obtained from the second order statistics of the processes  $x[n]$  and  $y[n]$ .

The problem of determining the estimator parameters by the LMMSE criterion is also called the Wiener filtering problem. Three subclasses of the problem are identified

1. The optimal smoothing problem  $N > 0$
2. The optimal filtering problem  $N = 0$
3. The optimal prediction problem  $N < 0$

In the smoothing problem, an estimate of the signal inside the duration of observation of the signal is made. The filtering problem estimates the current value of the signal on the basis of the present and past observations. The prediction problem addresses the issues of optimal prediction of the future value of the signal on the basis of present and past observations.

### 5.3 Wiener-Hopf Equations

The mean-square error of estimation is given by

$$\begin{aligned} Ee^2[n] &= E(x[n] - \hat{x}[n])^2 \\ &= E\left(x[n] - \sum_{i=-N}^{M-1} h[i]y[n-i]\right)^2 \end{aligned}$$

We have to minimize  $Ee^2[n]$  with respect to each  $h[i]$  to get the optimal estimation.

Corresponding minimization is given by

$$\frac{\partial E\{e^2[n]\}}{\partial h[j]} = 0, \text{ for } j = -N \dots 0 \dots M-1$$

( $E$  being a linear operator,  $E$  and  $\frac{\partial}{\partial h[j]}$  can be interchanged)

$$Ee[n]y[n-j] = 0, \quad j = -N \dots 0, 1, \dots M-1 \quad (1)$$

or

$$E \left( \overbrace{x[n] - \sum_{i=-N_a}^{M-1} h[i]y[n-i]}^{e[n]} \right) y[n-j] = 0, \quad j = -N \dots 0, 1, \dots M-1 \quad (2)$$

$$R_{XY}(j) = \sum_{i=-N_a}^{M-1} h[i]R_{YY}[j-i], \quad j = -N \dots 0, 1, \dots M-1 \quad (3)$$

This set of  $N+M+1$  equations in (3) are called Wiener Hopf equations or Normal equations.

- The result in (1) is the orthogonality principle which implies that the error is orthogonal to observed data.
- $\hat{x}[n]$  is the projection of  $x[n]$  onto the subspace spanned by observations  $y[n-M], y[n-M+1] \dots y[n], \dots y[n+N]$ .
- The estimation uses second order-statistics i.e. autocorrelation and cross-correlation functions.

- If  $x[n]$  and  $y[n]$  are jointly Gaussian then MMSE and LMMSE are equivalent. Otherwise we get a sub-optimum result.
- Also observe that

$$x[n] = \hat{x}[n] + e[n]$$

where  $\hat{x}[n]$  and  $e[n]$  are the parts of  $x[n]$  respectively correlated and uncorrelated with  $y[n]$ . Thus LMMSE separates out that part of  $x[n]$  which is correlated with  $y[n]$ . Hence the Wiener filter can be also interpreted as the correlation canceller. (See Orfanidis).

## 5.4 FIR Wiener Filter

$$\hat{x}[n] = \sum_{i=0}^{M-1} h[i]y[n-i]$$

The model parameters are given by the orthogonality principle

$$E \left( \overbrace{x[n] - \sum_{i=0}^{M-1} h[i]y[n-i]}^{e[n]} \right) y[n-j] = 0, \quad j = 0, 1, \dots, M-1$$

$$\sum_{i=0}^{M-1} h[i]R_{YY}[j-i] = R_{XY}(j), \quad j = 0, 1, \dots, M-1$$

In matrix form, we have

$$\mathbf{R}_{YY} \mathbf{h} = \mathbf{r}_{XY}$$

where

$$\mathbf{R}_{YY} = \begin{bmatrix} R_{YY}[0] & R_{YY}[-1] & \dots & R_{YY}[1-M] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[2-M] \\ \dots & \dots & \dots & \dots \\ R_{YY}[M-1] & R_{YY}[N-2] & \dots & R_{YY}[0] \end{bmatrix}$$

and

$$\mathbf{r}_{XY} = \begin{bmatrix} R_{XY}[0] \\ R_{XY}[1] \\ \dots \\ R_{XY}[M-1] \end{bmatrix}$$

and

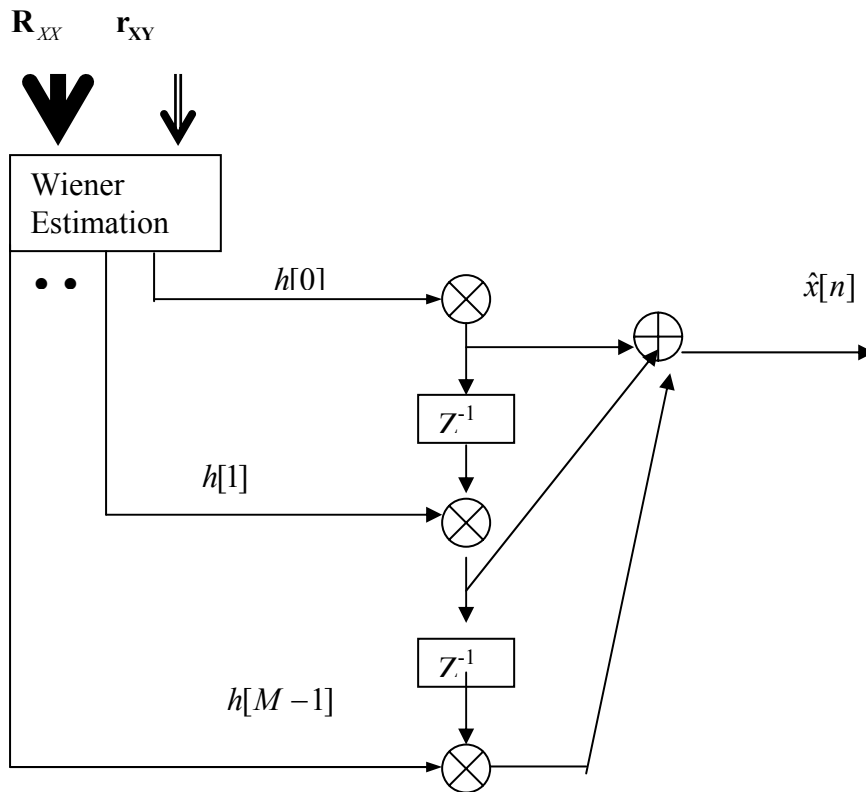
$$\mathbf{h} = \begin{bmatrix} h[0] \\ h[1] \\ \dots \\ h[M-1] \end{bmatrix}$$

Therefore,

$$\mathbf{h} = \mathbf{R}_{YY}^{-1} \mathbf{r}_{XY}$$

## 5.5 Minimum Mean Square Error - FIR Wiener Filter

$$\begin{aligned} E(e^2[n]) &= E\left\{e[n] \left( x[n] - \sum_{i=0}^{M-1} h[i]y[n-i] \right)\right\} \\ &= E\{e[n] x[n]\} \quad \because \text{error is orthogonal to data} \\ &= E\left\{ \left( x[n] - \sum_{i=0}^{M-1} h[i]y[n-i] \right) x[n] \right\} \\ &= R_{XX}[0] - \sum_{i=0}^{M-1} h[i]R_{XY}[i] \end{aligned}$$



### **Example1: Noise Filtering**

Consider the case of a carrier signal in presence of white Gaussian noise

$$x[n] = A \cos[w_0 n + \phi], \quad w_0 = \frac{\pi}{4}$$

$$y[n] = x[n] + v[n]$$

here  $\phi$  is uniformly distributed in  $(1, 2\pi)$ .

$v[n]$  is white Gaussian noise sequence of variance 1 and is independent of  $x[n]$ . Find the parameters for the FIR Wiener filter with  $M=3$ .

$$R_{xx}[m] = \frac{A^2}{2} \cos w_0 m$$

$$\begin{aligned} R_{yy}[m] &= E y[n]y[n-m] \\ &= E (x[n] + v[n])(x[n-m] + v[n-m]) \\ &= R_{xx}[m] + R_{vv}[m] + 0 + 0 \\ &= \frac{A^2}{2} \cos(w_0 m) + \delta[m] \end{aligned}$$

$$\begin{aligned} R_{xy}[m] &= E x[n]y[n-m] \\ &= E x[n](x[n-m] + v[n-m]) \\ &= R_{xx}[m] \end{aligned}$$

Hence the Wiener Hopf equations are

$$\begin{bmatrix} R_{yy}[0] & R_{yy}[1] & R_{yy}[2] \\ R_{yy}[1] & R_{yy}[0] & R_{yy}[1] \\ R_{yy}[2] & R_{yy}[1] & R_{yy}[0] \end{bmatrix} \begin{bmatrix} h[0] \\ h[1] \\ h[2] \end{bmatrix} = \begin{bmatrix} R_{xy}[0] \\ R_{xy}[1] \\ R_{xy}[2] \end{bmatrix}$$
$$\begin{bmatrix} \frac{A^2}{2} + 1 & \frac{A^2}{2} \cos \frac{\pi}{4} & \frac{A^2}{2} \cos \frac{\pi}{2} \\ \frac{A^2}{2} \cos \frac{\pi}{4} & \frac{A^2}{2} + 1 & \frac{A^2}{2} \cos \frac{\pi}{4} \\ \frac{A^2}{2} \cos \frac{\pi}{2} & \frac{A^2}{2} \cos \frac{\pi}{4} & \frac{A^2}{2} + 1 \end{bmatrix} \begin{bmatrix} h[0] \\ h[1] \\ h[2] \end{bmatrix} = \begin{bmatrix} \frac{A^2}{2} \\ \frac{A^2}{2} \cos \frac{\pi}{4} \\ \frac{A^2}{2} \cos \frac{\pi}{2} \end{bmatrix}$$

suppose  $A = 5\sqrt{2}$  then

$$\begin{bmatrix} 13.5 & \frac{12.5}{\sqrt{2}} & 0 \\ \frac{12.5}{\sqrt{2}} & 13.5 & \frac{12.5}{\sqrt{2}} \\ 0 & \frac{12.5}{\sqrt{2}} & 13.5 \end{bmatrix} \begin{bmatrix} h[0] \\ h[1] \\ h[2] \end{bmatrix} = \begin{bmatrix} 12.5 \\ \frac{12.5}{\sqrt{2}} \\ 0 \end{bmatrix}$$

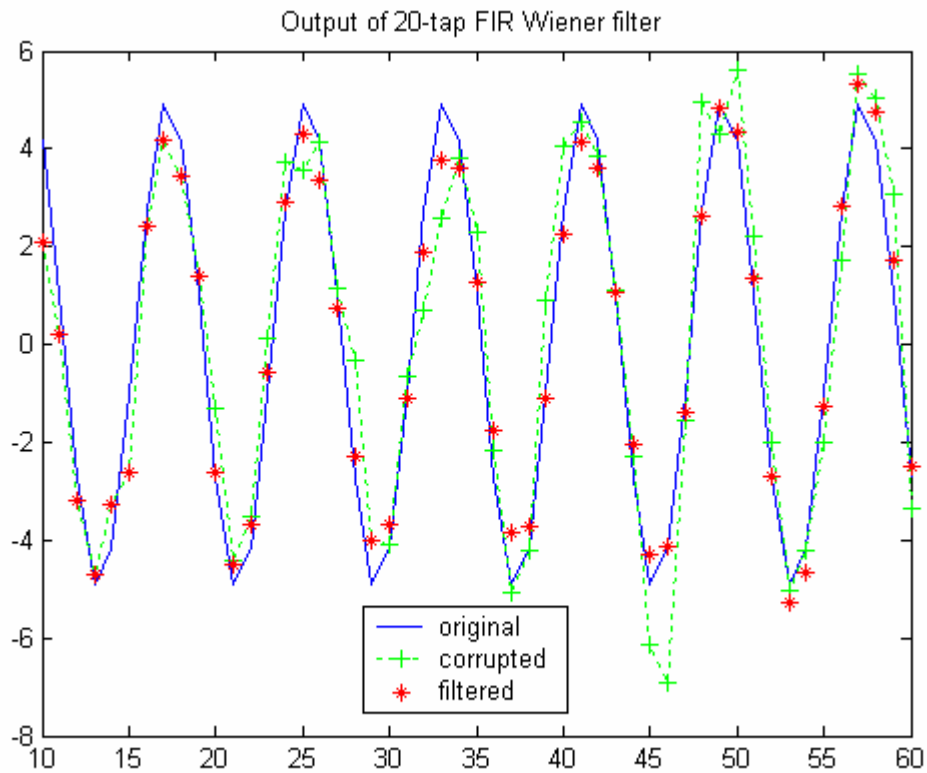
$$\begin{bmatrix} h[0] \\ h[1] \\ h[2] \end{bmatrix} = \begin{bmatrix} 13.5 & \frac{12.5}{\sqrt{2}} & 0 \\ \frac{12.5}{\sqrt{2}} & 13.5 & \frac{12.5}{\sqrt{2}} \\ 0 & \frac{12.5}{\sqrt{2}} & 13.5 \end{bmatrix}^{-1} \begin{bmatrix} 12.5 \\ 12.5 \\ 0 \end{bmatrix}$$

$$h[0] = 0.707$$

$$h[1] = 0.34$$

$$h[2] = -0.226$$

Plot the filter performance for the above values of  $h[0]$ ,  $h[1]$  and  $h[2]$ . The following figure shows the performance of the 20-tap FIR wiener filter for noise filtering.



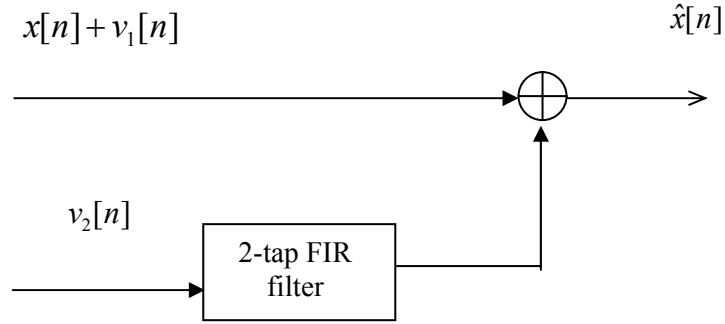


**Example 2 : Active Noise Control**

Suppose we have the observation signal  $y[n]$  is given by

$$y[n] = 0.5 \cos(w_0 n + \phi) + v_1[n]$$

where  $\phi$  is uniformly distributed in  $(0, 2\pi)$  and  $v_1[n] = 0.6v[n-1] + v[n]$  is an MA(1) noise. We want to control  $v_1[n]$  with the help of another correlated noise  $v_2[n]$  given by  $v_2[n] = 0.8v[n-1] + v[n]$



The Wiener Hopf Equations are given by

$$\mathbf{R}_{v_2 v_2} \mathbf{h} = \mathbf{r}_{v_1 v_2}$$

where  $\mathbf{h} = [h[0] \ h[1]]'$

and

$$\mathbf{R}_{v_2 v_2} = \begin{bmatrix} 1.64 & 0.8 \\ 0.8 & 1.64 \end{bmatrix} \text{ and}$$

$$\mathbf{r}_{v_1 v_2} = \begin{bmatrix} 1.48 \\ 0.6 \end{bmatrix}$$

$$\therefore \begin{bmatrix} h[0] \\ h[1] \end{bmatrix} = \begin{bmatrix} 0.9500 \\ -0.0976 \end{bmatrix}$$

**Example 3:**

(Continuous time prediction) Suppose we want to predict the continuous-time process

$X(t)$  at time  $(t + \tau)$  by

$$\hat{X}(t + \tau) = aX(t)$$

Then by orthogonality principle

$$E(X(t + \tau) - aX(t))X(t) = 0$$

$$\Rightarrow a = \frac{R_{XX}(\tau)}{R_{XX}(0)}$$

As a particular case consider the first-order Markov process given by

$$\frac{d}{dt} X(t) = AX(t) + v(t)$$

In this case,

$$R_{XX}(\tau) = R_{XX}(0)e^{-A\tau}$$

$$\therefore a = \frac{R_{XX}(\tau)}{R_{XX}(0)} = e^{-A\tau}$$

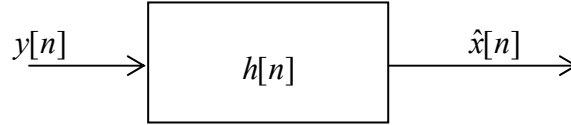
Observe that for such a process

$$\begin{aligned} E(X(t+\tau) - aX(t))X(t-\tau_1) &= 0 \\ &= R_{XX}(\tau+\tau_1) - aR_{XX}(\tau_1) \\ &= R_{XX}(0)e^{-A(\tau+\tau_1)} - e^{-A\tau}R_{XX}(0)e^{-A\tau_1} \\ &= 0 \end{aligned}$$

Therefore, the linear prediction of such a process based on any past value is same as the linear prediction based on current value.

## 5.6 IIR Wiener Filter (Causal)

Consider the IIR filter to estimate the signal  $x[n]$  shown in the figure below.



The estimator  $\hat{x}[n]$  is given by

$$\hat{x}(n) = \sum_{i=0}^{\infty} h(i)y(n-i)$$

The mean-square error of estimation is given by

$$\begin{aligned} Ee^2[n] &= E(x[n] - \hat{x}[n])^2 \\ &= E(x[n] - \sum_{i=0}^{\infty} h[i]y[n-i])^2 \end{aligned}$$

We have to minimize  $Ee^2[n]$  with respect to each  $h[i]$  to get the optimal estimation.

Applying the orthogonality principle, we get the WH equation.

$$E(x[n] - \sum_{i=0}^{\infty} h(i)y[n-i])y[n-j] = 0, \quad j = 0, 1, \dots$$

From which we get

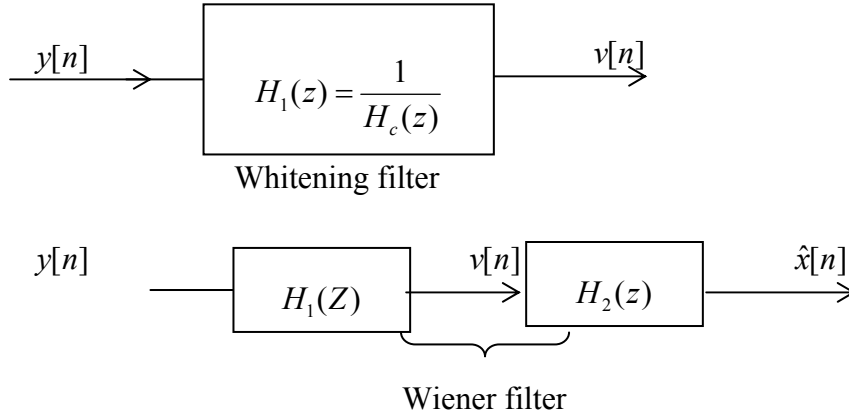
$$\sum_{i=0}^{\infty} h[i] R_{YY}[j-i] = R_{XY}[j], \quad j = 0, 1, \dots$$

- We have to find  $h[i]$ ,  $i = 0, 1, \dots, \infty$  by solving the above infinite set of equations.
- This problem is better solved in the z-transform domain, though we cannot directly apply the convolution theorem of z-transform.

Here comes Wiener's contribution.

The analysis is based on the spectral Factorization theorem:

$$S_{YY}(z) = \sigma_v^2 H_c(z) H_c(z^{-1})$$



Now  $h_2[n]$  is the coefficient of the Wiener filter to estimate  $x[n]$  from the innovation sequence  $v[n]$ . Applying the orthogonality principle results in the Wiener Hopf equation

$$\hat{x}(n) = \sum_{i=0}^{\infty} h_2(i) v(n-i)$$

$$E \left\{ x[n] - \sum_{i=0}^{\infty} h_2[i] v[n-i] \right\} v[n-j] = 0$$

$$\therefore \sum_{i=0}^{\infty} h_2[i] R_{VV}[j-i] = R_{XV}[j], \quad j = 0, 1, \dots$$

$$R_{VV}[m] = \sigma_v^2 \delta[m]$$

$$\therefore \sum_{i=0}^{\infty} h_2(i) \sigma_v^2 \delta[j-i] = R_{XV}(j), \quad j = 0, 1, \dots$$

So that

$$h_2[j] = \frac{R_{XV}[j]}{\sigma_v^2} \quad j \geq 0$$

$$H_2(z) = \frac{[S_{XV}(z)]_+}{\sigma_v^2}$$

where  $[S_{XV}(z)]_+$  is the positive part (i.e., containing non-positive powers of  $z$ ) in power series expansion of  $S_{XV}(z)$ .

$$v[n] = \sum_{i=0}^{\infty} h_1[i] y[n-i]$$

$$\begin{aligned} R_{xv}[j] &= Ex[n]v[n-j] \\ &= \sum_{i=0}^{\infty} h_1[i] E x[n] y[n-j-i] \\ &= \sum_{i=0}^{\infty} h_1[i] R_{xy}[j+i] \end{aligned}$$

$$S_{xv}(z) = H_1(z^{-1}) S_{xy}(z) = \frac{1}{H_c(z^{-1})} S_{xy}(z)$$

$$\therefore H_2(z) = \frac{1}{\sigma_v^2} \left[ \frac{S_{xy}(z)}{H_c(z^{-1})} \right]_+$$

Therefore,

$$H(z) = H_1(z) H_2(z) = \frac{1}{\sigma_v^2 H_c(z)} \left[ \frac{S_{xy}(z)}{H_c(z^{-1})} \right]_+$$

We have to

- find the power spectrum of data and the cross power spectrum of the of the desired signal and data from the available model or estimate them from the data
- factorize the power spectrum of the data using the spectral factorization theorem

## 5.7 Mean Square Estimation Error – IIR Filter (Causal)

$$\begin{aligned} E(e^2[n]) &= E e[n] \left( x[n] - \sum_{i=0}^{\infty} h[i] y[n-i] \right) \\ &= E e[n] x[n] \quad \because \text{error is orthogonal to data} \\ &= E \left( x[n] - \sum_{i=0}^{\infty} h[i] y[n-i] \right) x[n] \\ &= R_{xx}[0] - \sum_{i=0}^{\infty} h[i] R_{xy}[i] \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} S_x(w) dw - \frac{1}{2\pi} \int_{-\pi}^{\pi} H(w) S_{xy}^*(w) dw \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} (S_x(w) - H(w) S_{xy}^*(w)) dw \\ &= \frac{1}{2\pi} \oint_C (S_x(z) - H(z) S_{xy}(z^{-1})) z^{-1} dz \end{aligned}$$

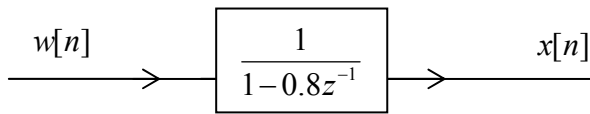
**Example 4:**

$y[n] = x[n] + v_1[n]$  observation model with

$$x[n] = 0.8 x[n-1] + w[n]$$

where  $v_1[n]$  is an additive zero-mean Gaussian white noise with variance 1 and  $w[n]$  is zero-mean white noise with variance 0.68. Signal and noise are uncorrelated.

Find the optimal Causal Wiener filter to estimate  $x[n]$ .

**Solution:**

$$S_{xx}(z) = \frac{0.68}{(1 - 0.8z^{-1})(1 - 0.8z)}$$

$$\begin{aligned} R_{yy}[m] &= E[y[n]y[n-m]] \\ &= E[(x[n] + v_1[n])(x[n-m] + v_1[n-m])] \\ &= R_{xx}[m] + R_{v_1v_1}[m] + 0 + 0 \end{aligned}$$

$$S_{yy}(z) = S_{xx}(z) + 1$$

Factorize

$$\begin{aligned} S_{yy}(z) &= \frac{0.68}{(1 - 0.8z^{-1})(1 - 0.8z)} + 1 \\ &= \frac{2(1 - 0.4z^{-1})(1 - 0.4z)}{(1 - 0.8z^{-1})(1 - 0.8z)} \\ \therefore H_c(z) &= \frac{(1 - 0.4z^{-1})}{(1 - 0.8z^{-1})} \end{aligned}$$

and

$$\sigma_v^2 = 2$$

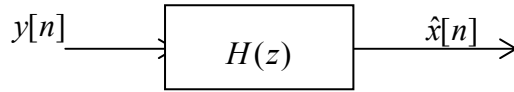
Also

$$\begin{aligned} R_{xy}[m] &= E[x[n]y[n-m]] \\ &= E[x[n](x[n-m] + v_1[n-m])] \\ &= R_{xx}[m] \end{aligned}$$

$$\begin{aligned} S_{xy}(z) &= S_{xx}(z) \\ &= \frac{0.68}{(1 - 0.8z^{-1})(1 - 0.8z)} \end{aligned}$$

$$\begin{aligned}
\therefore H(z) &= \frac{1}{\sigma_v^2 H_c(z)} \left[ \frac{S_{XY}(z)}{H_c(z^{-1})} \right]_+ \\
&= \frac{1}{2} \frac{(1-0.8z^{-1})}{(1-0.4z^{-1})} \left[ \frac{0.68}{(1-0.8z^{-1})(1-0.8z)} \right]_+ \\
&= \frac{0.944}{(1-0.4z^{-1})} \\
h[n] &= 0.944(0.4)^n \quad n \geq 0
\end{aligned}$$

## 5.8 IIR Wiener filter (Noncausal)



The estimator  $\hat{x}[n]$  is given by

$$\hat{x}[n] = \sum_{i=-\alpha}^{\alpha} h[i] y[n-i]$$

For LMMSE, the error is orthogonal to data.

$$E \left( x[n] - \sum_{i=-\alpha}^{\alpha} h[i] y[n-i] \right) y[n-j] = 0 \quad \forall j \in I$$

$$\sum_{i=-\infty}^{\infty} h[i] R_{YY}[j-i] = R_{XY}[j], \quad j = -\infty, \dots, 0, 1, \dots, \infty$$

- This form Wiener Hopf Equation is simple to analyse.
- Easily solved in frequency domain. So taking Z transform we get
- Not realizable in real time

$$H(z)S_{YY}(z) = S_{XY}(z)$$

so that

$$H(z) = \frac{S_{XY}(z)}{S_{YY}(z)}$$

or

$$H(w) = \frac{S_{XY}(w)}{S_{YY}(w)}$$

## 5.9 Mean Square Estimation Error – IIR Filter (Noncausal)

The mean square error of estimation is given by

$$\begin{aligned}
 E(e^2[n]) &= Ee[n] \left( x[n] - \sum_{i=-\infty}^{\infty} h[i]y[n-i] \right) \\
 &= Ee[n] x[n] \quad \because \text{error is orthogonal to data} \\
 &= E \left( x[n] - \sum_{i=-\infty}^{\infty} h[i]y[n-i] \right) x[n] \\
 &= R_{XX}[0] - \sum_{i=-\infty}^{\infty} h[i]R_{XY}[i] \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} S_X(w) dw - \frac{1}{2\pi} \int_{-\pi}^{\pi} H(w) S_{XY}^*(w) dw \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} (S_X(w) - H(w) S_{XY}^*(w)) dw \\
 &= \frac{1}{2\pi} \oint_C (S_X(z) - H(z) S_{XY}(z^{-1})) z^{-1} dz
 \end{aligned}$$

### **Example 5: Noise filtering by noncausal IIR Wiener Filter**

Consider the case of a carrier signal in presence of white Gaussian noise

$$y[n] = x[n] + v[n]$$

where  $v[n]$  is additive zero-mean Gaussian white noise with variance  $\sigma_v^2$ . Signal and noise are uncorrelated

$$S_{YY}(w) = S_{XX}(w) + S_{VV}(w)$$

and

$$S_{XY}(w) = S_{XX}(w)$$

$$\begin{aligned}
 \therefore H(w) &= \frac{S_{XX}(w)}{S_{XX}(w) + S_{VV}(w)} \\
 &= \frac{\frac{S_{XY}(w)}{S_{VV}(w)}}{\frac{S_{XX}(w)}{S_{VV}(w)} + 1}
 \end{aligned}$$

Suppose SNR is very high

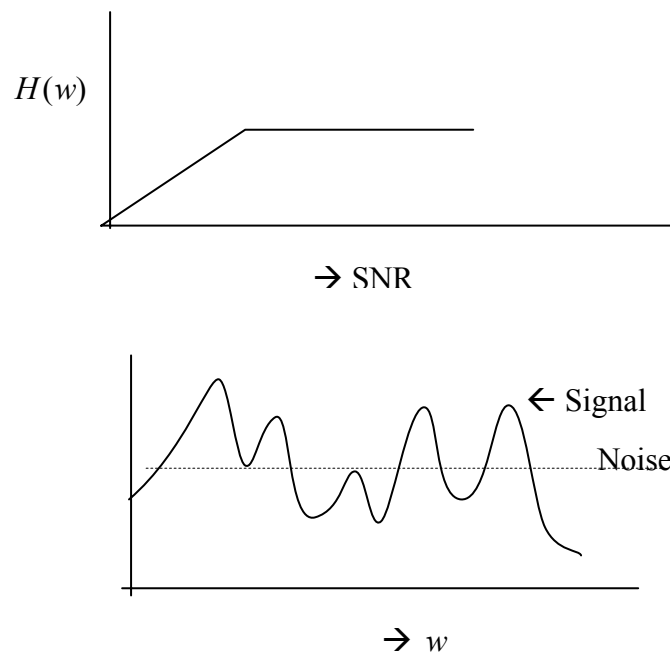
$$H(w) \cong 1$$

(i.e. the signal will be passed un-attenuated).

When SNR is low

$$H(w) = \frac{S_{XX}(w)}{S_{VV}(w)}$$

(i.e. If noise is high the corresponding signal component will be attenuated in proportion of the estimated SNR.



**Figure -** (a) a high-SNR signal is passed unattended by the IIR Wiener filter  
(b) Variation of SNR with frequency

### **Example 6: Image filtering by IIR Wiener filter**

$S_{YY}(w)$  = power spectrum of the corrupted image

$S_{VV}(w)$  = power spectrum of the noise, estimated from the noise model  
or from the constant intensity ( like back-ground) of the image

$$H(w) = \frac{S_{XX}(w)}{S_{XX}(w) + S_{VV}(w)}$$

$$= \frac{S_{YY}(w) - S_{VV}(w)}{S_{YY}(w)}$$

### **Example 7:**

Consider the signal in presence of white noise given by

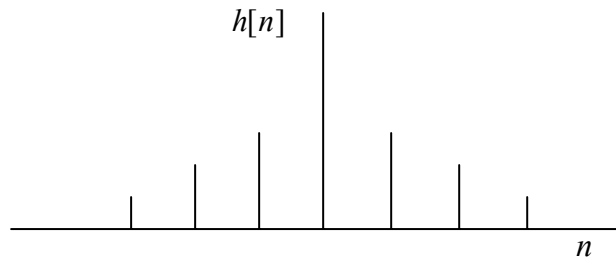
$$x[n] = 0.8 x[n-1] + w[n]$$

where  $v[n]$  is an additive zero-mean Gaussian white noise with variance 1 and  $w[n]$  is zero-mean white noise with variance 0.68. Signal and noise are uncorrelated.

Find the optimal noncausal Wiener filter to estimate  $x[n]$ .



$$\begin{aligned}
 H(z) &= \frac{S_{XY}(z)}{S_{YY}(z)} = \frac{\frac{0.68}{(1-0.8z^{-1})(1-0.8z)}}{2 \frac{(1-0.4z^{-1})(1-0.4z)}{(1-0.6z^{-1})(1-0.6z)}} \\
 &= \frac{0.34}{(1-0.4z^{-1})(1-0.4z)} \quad \text{One pole outside the unit circle} \\
 &= \frac{0.4048}{1-0.4z^{-1}} + \frac{0.4048}{1-0.4z} \\
 \therefore h[n] &= 0.4048(0.4)^n u(n) + 0.4048(0.4)^{-n} u(-n-1)
 \end{aligned}$$



**Figure - Filter Impulse Response**

## **CHAPTER – 6: LINEAR PREDICTION OF SIGNAL**

### **6.1 Introduction**

Given a sequence of observation

$y[n-1], y[n-2], \dots, y[n-M]$ , what is the best prediction for  $y[n]$ ?

(one-step ahead prediction)

The minimum mean square error prediction  $\hat{y}[n]$  for  $y[n]$  is given by  $\hat{y}[n] = E \{y[n] | y[n-1], y[n-2], \dots, y[n-M]\}$  which is a nonlinear predictor.

A linear prediction is given by

$$\hat{y}[n] = \sum_{i=1}^M h[i] y[n-i]$$

where  $h[i], i = 1 \dots M$  are the prediction parameters.

- Linear prediction has very wide range of applications.
- For an exact AR (M) process, linear prediction model of order M and the corresponding AR model have the same parameters. For other signals LP model gives an approximation.

### **6.2 Areas of application**

- Speech modeling
- ECG modeling
- Low-bit rate speech coding
- DPCM coding
- Speech recognition
- Internet traffic prediction

LPC (10) is the popular linear prediction model used for speech coding. For a frame of speech samples, the prediction parameters are estimated and coded. In CELP (Code book Excited Linear Prediction) the prediction error  $e[n] = y[n] - \hat{y}[n]$  is vector quantized and transmitted.

$\hat{y}[n] = \sum_{i=1}^M h[i] y[n-i]$  is a FIR Wiener filter shown in the following figure. It is called the *linear prediction filter*.

Therefore

$$\begin{aligned} e[n] &= y[n] - \hat{y}[n] \\ &= y[n] - \sum_{i=1}^M h[i] y[n-i] \end{aligned}$$

is the prediction error and the corresponding filter is called prediction error filter.

Linear Minimum Mean Square error estimates for the prediction parameters are given by the orthogonality relation

$$E\{e[n] y[n-j]\} = 0 \quad \text{for } j = 1, 2, \dots, M$$

$$\therefore E\left(y[n] - \sum_{i=1}^M h[i] y[n-i]\right) y[n-j] = 0 \quad j = 1, 2, \dots, M$$

$$\Rightarrow R_{yy}[j] - \sum_{i=1}^M h[i] R_{yy}[j-i] = 0$$

$$\Rightarrow R_{yy}[j] = \sum_{i=1}^M h[i] R_{yy}[j-i] \quad j = 1, 2, \dots, M$$

which is the Wiener Hopf equation for the linear prediction problem and same as the Yule Walker equation for AR (M) Process.

In Matrix notation

$$\begin{bmatrix} R_{yy}[0] & R_{yy}[1] & \dots & R_{yy}[M-1] \\ R_{yy}[1] & R_{yy}[0] & \dots & R_{yy}[M-2] \\ \vdots & \vdots & \ddots & \vdots \\ R_{yy}[M-1] & R_{yy}[M-2] & \dots & R_{yy}[0] \end{bmatrix} \begin{bmatrix} h[1] \\ h[2] \\ \vdots \\ h[M] \end{bmatrix} = \begin{bmatrix} R_{yy}[1] \\ R_{yy}[2] \\ \vdots \\ R_{yy}[M] \end{bmatrix}$$

$$\mathbf{R}_{yy} \mathbf{h} = \mathbf{r}_{yy}$$

$$\therefore \mathbf{h} = (\mathbf{R}_{yy})^{-1} \mathbf{r}_{yy}$$

### 6.3 Mean Square Prediction Error (MSPE)

$$\begin{aligned} E(e^2[n]) &= E\left(y[n] - \sum_{i=1}^M h[i] y[n-i]\right) e[n] \\ &= E\{y[n] e[n]\} \\ &= E\left\{y[n] \left(y[n] - \sum_{i=1}^M h[i] y[n-i]\right)\right\} \\ &= R_{yy}[0] - \sum_{i=1}^M h[i] R_{yy}[i] \end{aligned}$$

## 6.4 Forward Prediction Problem

The above linear prediction problem is the forward prediction problem. For notational simplicity let us rewrite the prediction equation as

$$\hat{y}[n] = \sum_{i=1}^M h_M[i] y[n-i]$$

where the prediction parameters are being denoted by  $h_M[i], i = 1 \dots M$ .

## 6.5 Backward Prediction Problem

Given  $y[n], y[n-1], \dots, y[n-M+1]$ , we want to estimate  $y[n-M]$ .

The linear prediction is given by

$$\hat{y}[n-M] = \sum_{i=1}^M b_M[i] y[n+1-i]$$

Applying orthogonality principle.

$$E(y[n-M] - \sum_{i=1}^M b_M[i] y[n+1-i]) y[n+1-j] = 0 \quad j = 1, 2, \dots, M.$$

This will give

$$R_{YY}[M+1-j] = \sum_{i=1}^M b_M[i] R_{YY}[j-i] \quad j = 1, 2, \dots, M$$

Corresponding matrix form

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] & \dots & R_{YY}[M-1] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[M-2] \\ \cdot & \cdot & \cdot & \cdot \\ R_{YY}[M-1] & R_{YY}[M-2] & \dots & R_{YY}[0] \end{bmatrix} \begin{bmatrix} b_M[1] \\ b_M[2] \\ \cdot \\ b_M[M] \end{bmatrix} = \begin{bmatrix} R_{YY}[M] \\ R_{YY}[M-1] \\ \cdot \\ R_{YY}[1] \end{bmatrix} \quad (1)$$

## 6.6 Forward Prediction

Rewriting the Mth-order forward prediction problem, we have

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] & \dots & R_{YY}[M-1] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[M-2] \\ \cdot & \cdot & \cdot & \cdot \\ R_{YY}[M-1] & R_{YY}[M-2] & \dots & R_{YY}[0] \end{bmatrix} \begin{bmatrix} h_M[1] \\ h_M[2] \\ \cdot \\ h_M[M] \end{bmatrix} = \begin{bmatrix} R_{YY}[1] \\ R_{YY}[2] \\ \cdot \\ R_{YY}[M] \end{bmatrix} \quad (2)$$

From (1) and (2) we conclude

$$b_M[i] = h_M[M+1-i], i = 1, 2, \dots, M$$

Thus forward prediction parameters in reverse order will give the backward prediction parameters.

M S prediction error

$$\begin{aligned}\varepsilon_M &= E \left( y[n-M] - \sum_{i=1}^M b_M[i] y[n+1-i] \right) y[n-M] \\ &= R_{YY}[0] - \sum_{i=1}^M b_M[i] R_{YY}[M+1-i] \\ &= R_{YY}[0] - \sum_{i=1}^M h_M[M+1-i] R_{YY}[M+1-i]\end{aligned}$$

which is same as the forward prediction error.

Thus

Backward prediction error = Forward Prediction error.
---

### **Example 1:**

Find the second order predictor for  $y[n]$  given  $y[n] = x[n] + v[n]$ , where  $v[n]$  is a 0-mean white noise with variance 1 and uncorrelated with  $x[n]$  and  $x[n] = 0.8x[n-1] + w[n]$ ,  $w[n]$  is a 0-mean random variable with variance 0.68

The linear predictor is given by

$$\hat{y}[n] = h_2[1] y[n-1] + h_2[2] y[n-2]$$

We have to find  $h_2[1]$  and  $h_2[2]$ .

Corresponding Yule Walker equations are

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] \\ R_{YY}[1] & R_{YY}[0] \end{bmatrix} \begin{bmatrix} h_2[1] \\ h_2[2] \end{bmatrix} = \begin{bmatrix} R_{YY}[1] \\ R_{YY}[2] \end{bmatrix}$$

To find out  $R_{YY}[0]$ ,  $R_{YY}[1]$  and  $R_{YY}[2]$

$$y[n] = x[n] + v[n],$$

$$R_{YY}[m] = R_{XX}[m] + \delta[m]$$

$$x[n] = 0.8x[n-1] + w[n]$$

$$\therefore R_{XX}[m] = \frac{0.68}{1 - (0.8)^2} (0.8)^{|m|} = 1.89 \times (0.8)^{|m|}$$

$$R_{YY}[0] = 2.89, R_{YY}[1] = 1.51 \text{ and } R_{YY}[2] = 1.21$$

$$\text{Solving } h_2[1] = 0.4178 \text{ and } h_2[2] = 0.2004.$$

## 6.7 Levinson Durbin Algorithm

Levinson Durbin algorithm is the most popular technique for determining the LPC parameters from a given autocorrelation sequence.

Consider the Yule Walker equation for  $m$ th order linear predictor.

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] & \dots & R_{YY}[m-1] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[m-2] \\ \cdot & \cdot & \cdot & \cdot \\ R_{YY}[m-1] & \dots & \dots & R_{YY}[0] \end{bmatrix} \begin{bmatrix} h_m[1] \\ h_m[2] \\ \cdot \\ h_m[m] \end{bmatrix} = \begin{bmatrix} R_{YY}[1] \\ \cdot \\ \cdot \\ R_{YY}[m] \end{bmatrix} \quad (1)$$

Writing in the reverse order

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] & \dots & R_{YY}[m-1] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[m-2] \\ \cdot & \cdot & \cdot & \cdot \\ R_{YY}[m-1] & \dots & \dots & R_{YY}[0] \end{bmatrix} \begin{bmatrix} h_m[m] \\ h_m[m-1] \\ \cdot \\ h_m[1] \end{bmatrix} = \begin{bmatrix} R_{YY}[m] \\ \cdot \\ \cdot \\ R_{YY}[1] \end{bmatrix} \quad (2)$$

Then  $(m+1)$  the order predictor is given by

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] & \dots & R_{YY}[m-1] & R_{YY}[m] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[m-2] & R_{YY}[m-1] \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ R_{YY}[m-1] & R_{YY}[m-2] & \dots & R_{YY}[0] & R_{YY}[1] \\ R_{YY}[m] & R_{YY}[m-1] & \dots & R_{YY}[1] & R_{YY}[0] \end{bmatrix} \begin{bmatrix} h_{m+1}[1] \\ h_{m+1}[2] \\ \cdot \\ h_{m+1}[m] \\ h_{m+1}[m+1] \end{bmatrix} = \begin{bmatrix} R_{YY}[1] \\ R_{YY}[2] \\ \cdot \\ R_{YY}[m] \\ R_{YY}[m+1] \end{bmatrix}$$

Let us partition equation (2) as shown. Then

$$\begin{bmatrix} R_{YY}[0] & R_{YY}[1] & \dots & R_{YY}[m-1] \\ R_{YY}[1] & R_{YY}[0] & \dots & R_{YY}[m-2] \\ \cdot & \cdot & \cdot & \cdot \\ R_{YY}[m-1] & R_{YY}[m-2] & \dots & R_{YY}[0] \end{bmatrix} \begin{bmatrix} h_{m+1}[1] \\ h_{m+1}[2] \\ \cdot \\ h_{m+1}[m] \end{bmatrix} + h_{m+1}[m+1] \begin{bmatrix} R_{YY}[m] \\ R_{YY}[m-1] \\ \cdot \\ R_{YY}[1] \end{bmatrix} = \begin{bmatrix} R_{YY}[1] \\ R_{YY}[2] \\ \cdot \\ R_{YY}[m] \end{bmatrix} \quad (3)$$

and

$$\sum_{i=1}^m h_{m+1}[i] R_{YY}[m+1-i] + h_{m+1}[m+1] R_{YY}[0] = R_{YY}[m+1] \quad (4)$$

From equation (3) premultiplying by  $\mathbf{R}_{YY}^{-1}$ , we get

$$\begin{bmatrix} h_{m+1}[1] \\ h_{m+1}[2] \\ \vdots \\ h_{m+1}[m] \end{bmatrix} + h_{m+1}[m+1] \mathbf{R}_{YY}^{-1} \begin{bmatrix} R_{YY}[m] \\ R_{YY}[m-1] \\ \vdots \\ R_{YY}[1] \end{bmatrix} = \mathbf{R}_{YY}^{-1} \begin{bmatrix} R_{YY}[1] \\ R_{YY}[2] \\ \vdots \\ R_{YY}[m] \end{bmatrix}$$

$$\begin{bmatrix} h_{m+1}[1] \\ h_{m+1}[2] \\ \vdots \\ h_{m+1}[m] \end{bmatrix} + h_{m+1}[m+1] \begin{bmatrix} h_m[m] \\ h_m[m-1] \\ \vdots \\ h_m[1] \end{bmatrix} = \begin{bmatrix} h_m[1] \\ h_m[2] \\ \vdots \\ h_m[m] \end{bmatrix}$$

The equations can be rewritten as

$$h_{m+1}[i] = h_m[i] + k_{m+1} h_m[m+1-i] \quad i = 1, 2, \dots, m \quad (5)$$

where  $k_m = -h_m[m]$  is called the reflection coefficient or the PARCOR (partial correlation) coefficient.

From equation (4) we get

$$\sum_{i=1}^m h_{m+1}[i] R_{YY}[m+1-i] + h_{m+1}[m] R_{YY}[0] = R_{YY}[m+1]$$

using equation (5)

$$\begin{aligned} \sum_{i=1}^m \{h_m[i] + k_{m+1} h_m[m+1-i]\} R_{YY}[m+1-i] - k_{m+1} R_{YY}[0] &= R_{YY}[m+1] \\ \sum_{i=1}^m h_m[i] R_{YY}[m+1-i] - k_{m+1} R_{YY}[0] + k_{m+1} \sum_{i=1}^m h_m[m+1-i] R_{YY}[m+1-i] &= R_{YY}[m+1] \\ k_{m+1} \{R_{YY}[0] - \sum_{i=1}^m h_m[m+1-i] R_{YY}[m+1-i]\} &= -R_{YY}[m+1] + \sum_{i=1}^m h_m[i] R_{YY}[m+1-i] \\ k_{m+1} &= \frac{-R_{YY}[m+1] + \sum_{i=1}^m h_m[i] R_{YY}[m+1-i]}{\mathcal{E}[m]} \\ &= \frac{\sum_{i=0}^m h_m[i] R_{YY}[m+1-i]}{\mathcal{E}[m]} \end{aligned}$$

where

$$\mathcal{E}[m] = R_{YY}[0] - \sum_{i=1}^m h_m[m+1-i] R_{YY}[m+1-i] \text{ is the mean - square prediction error.}$$

Here we have used the assumption that  $h_m[0] = -1$

$$\therefore \varepsilon[m+1] = R_{YY}[0] - \sum_{i=1}^{m+1} h_{m+1}[m+2-i] R_{YY}[m+2-i]$$

Using the recursion for  $h_{m+1}[i]$

We get

$$\varepsilon[m+1] = \varepsilon[m](1 - k_{m+1}^2)$$

will give MSE recursively. Since MSE is non negative

$$\begin{aligned} k_m^2 &\leq 1 \\ \therefore |k_m| &\leq 1 \end{aligned}$$

- If  $|k_m| < 1$ , the LPC error filter will be minimum-phase, and hence the corresponding synthesis filter will be stable.
- Efficient realization can be achieved in terms of  $k_m$ .
- $k_m$  represents the direct correlation of the data  $y[n-m]$  on  $y[n]$  when the correlation due to the intermediate data  $y[n-m+1], y[n-m+2], \dots, y[n-1]$  is removed. It is defined by

$$k_m = \frac{E e_m^f[n] e_m^f[n]}{R_{yy}(0)}$$

where

$$e_m^f[n] = \text{forward prediction error} = y[n] - \sum_{i=1}^n h_m[i] y[n-i]$$

and

$$e_m^b[n] = \text{backward prediction error} = y[n-m] - \sum_{i=1}^n h_m[m+1-i] y[n+1-i]$$

## 6.8 Steps of the Levinson- Durbin algorithm

Given  $R_{YY}[m], m = 0, 1, 2, \dots$

Initialization

Take  $h_m[0] = -1$  for all  $m$

For  $m = 0$ ,

$$\varepsilon[0] = R_{YY}[0]$$

For  $m = 1, 2, 3, \dots$

$$k_m = \frac{\sum_{i=0}^{m-1} h_{m-1}[i] R_{YY}[m-i]}{\varepsilon[m-1]}$$



$$h_m[i] = h_{m-1}[i] + k_m h_{m-1}[m-i], i = 1, 2, \dots, m-1$$

$$h_m[m] = -k_m$$

$$\varepsilon_m = \varepsilon_{m-1}(1 - k_m^2)$$

Go on computing up to given final value of  $m$ .

### **Some salient points**

- The reflection parameters and the mean-square error completely determine the LPC coefficients. Alternately, given the reflection coefficients and the final mean-square prediction error, we can determine the LPC coefficients.
- The algorithm is order recursive. By solving for  $m$ -th order linear prediction problem we get all previous order solutions
- If the estimated autocorrelation sequence satisfy the properties of an autocorrelation functions, the algorithm will yield stable coefficients

## **6.9 Lattice filter realization of Linear prediction error filters**

$e_m^f[n]$  = prediction error due to  $m$ th order forward prediction.

$e_m^b[n]$  = prediction error due to  $m$ th order backward prediction.

Then,

$$e_m^f[n] = y[n] - \sum_{i=1}^m h_m[i] y[n-i] \quad (1)$$

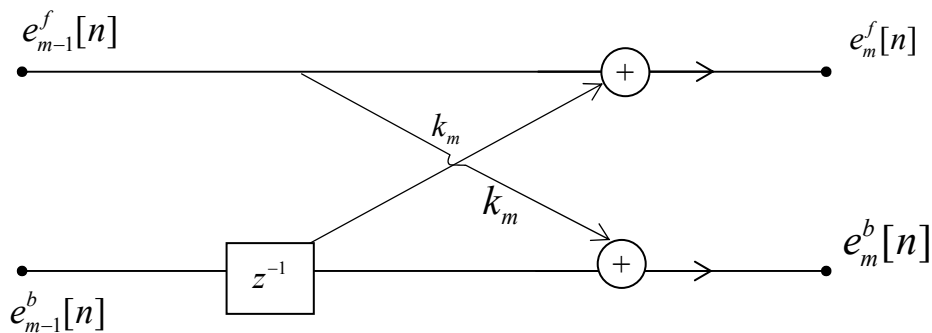
$$e_m^b[n] = y[n-m] - \sum_{i=1}^m h_m[m+1-i] y[n+1-i]$$

From (1), we get

$$\begin{aligned} e_m^f[n] &= y[n] - h_m[m]y[n-m] - \sum_{i=1}^{m-1} h_m[i]y[n-i] \\ &= y[n] + k_m y[n-m] - \sum_{i=1}^{m-1} (h_{m-1}[i] + k_m h_{m-1}[m-i])y[n-i] \\ &= y[n] - \sum_{i=1}^{m-1} h_{m-1}[i]y[n-i] + k_m \left( y[n-m] - \sum_{i=1}^{m-1} h_{m-1}[m-i]y[n-i] \right) \\ &= e_{m-1}^f[n] + k_m e_{m-1}^b[n-1] \\ \therefore e_m^f[n] &= e_{m-1}^f[n] + k_m e_{m-1}^b[n-1] \end{aligned}$$

Similarly we can show that

$$e_m^b[n] = e_{m-1}^b[n-1] + k_m e_{m-1}^f[n]$$



### How to initialize the lattice?

We have

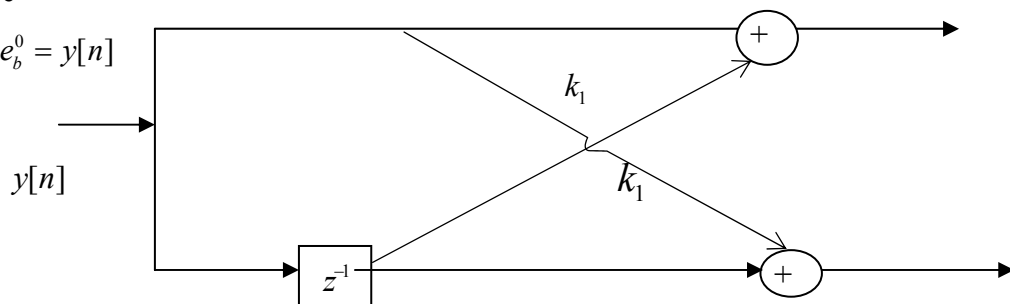
$$e_0^f = y[n] - 0$$

and

$$e_0^b = y[n] - 0$$

Hence

$$e_0^f = e_0^b = y[n]$$



## 6.10 Advantage of Lattice Structure

- Modular structure can be extended by first cascading another section. New stages can be added without modifying the earlier stages.
- Same elements are used in each stage. So efficient for VLSI implementation.
- Numerically efficient as  $|k_m| < 1$ .
- Each stage is decoupled with earlier stages.

=> It follows from the fact that for W.S.S. signal,  $e_m[n]$  sequences as a function of  $m$  are uncorrelated.

$$e_i^b[n] \text{ and } e_m^b[n] \quad 0 \leq i < m$$

are uncorrelated (Gram Schmidt orthogonalisation may be obtained through Lattice filter).

$$e_k^b[n] = y[n-k] - \sum_{i=1}^k h_k[k+1-i] y[n+1-i]$$

$$\begin{aligned} E e_m^b[n] e_k^b[n] &= E e_m^b[n] (y[n-k] - \sum_{i=1}^k h_k[k+1-i] y[n+1-i]) \\ &= 0 \quad \text{for } 0 \leq k < m \end{aligned}$$

Thus the lattice filter can be used to whiten a sequence.

With this result, it can be shown that

$$k_m = -\frac{E(e_{m-1}^f[n] e_{m-1}^b[n-1])}{E(e_{m-1}^b[n-1])^2} \quad (i)$$

and

$$k_m = -\frac{E(e_{m-1}^f[n] e_{m-1}^b[n-1])}{E(e_{m-1}^f[n])^2} \quad (ii)$$

### **Proof:**

Mean Square Prediction Error

$$\begin{aligned} &= E(e_m^f[n])^2 \\ &= E(e_{m-1}^f[n] + k_m e_{m-1}^b[n-1])^2 \end{aligned}$$

This is to be minimized w.r.t.  $k_m$

$$\begin{aligned} &\Rightarrow 2(e_{m-1}^f[n] + k_m e_{m-1}^b[n-1]) e_{m-1}^b[n-1] = 0 \\ &\Rightarrow (i) \end{aligned}$$

$$\text{Minimizing } E(e_m^b[n])^2 \Rightarrow (ii)$$

### **Example 2:**

Consider the random signal model  $y[n] = x[n] + v[n]$ , where  $v[n]$  is a 0-mean white noise with variance 1 and uncorrelated with  $x[n]$  and  $x[n] = 0.8x[n-1] + w[n]$ ,  $w[n]$  is a 0-mean random variable with variance 0.68

- Find the second –order linear predictor for  $y[n]$
- Obtain the lattice structure for the prediction error filter
- Use the above structure to design a second-order FIR Wiener filter to estimate  $x[n]$  from  $y[n]$ .

## CHAPTER – 7: ADAPTIVE FILTERS

### 7.1 Introduction

In practical situations, the system is operating in an uncertain environment where the input condition is not clear and/or the unexpected noise exists. Under such circumstances, the system should have the flexible ability to modify the system parameters and makes the adjustments based on the input signal and the other relevant signal to obtain optimal performance.

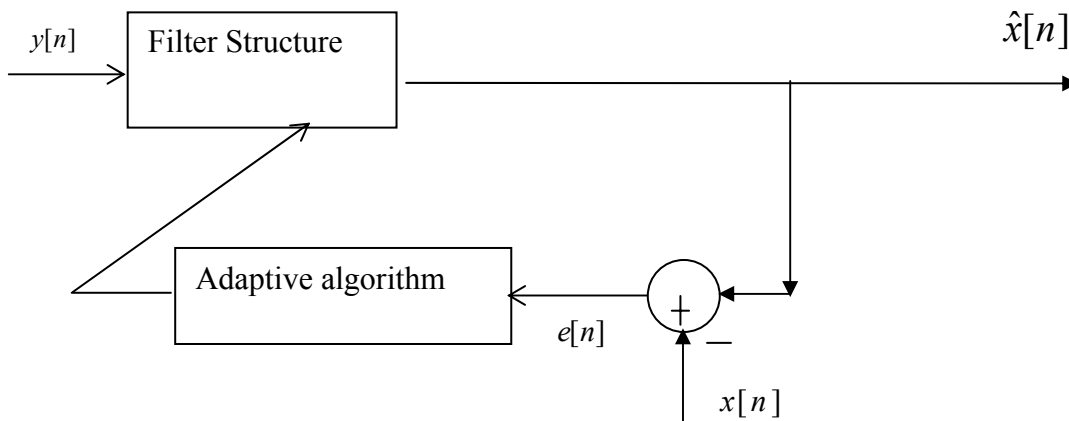
A system that searches for improved performance guided by a computational algorithm for adjustment of the parameters or weights is called an *adaptive system*. The adaptive system is time-varying.

Wiener filter is a linear time-invariant filter.

In practical situation, the signal is non-stationary. Under such circumstances, optimal filter should be time varying. The filter should have the ability to modify its parameters based on the input signal and the other relevant signal to obtain optimal performance.

How to do this?

- Assume stationarity within certain data length. Buffering of data is required and may work in some applications.
- The time-duration over which stationarity is a valid assumption, may be short so that accurate estimation of the model parameters is difficult.
- One solution is *adaptive filtering*. Here the filter coefficients are updated as a function of the filtering error. The basic filter structure is as shown in Fig. 1.



The filter structure is FIR of known tap-length, because the adaptation algorithm updates each filter coefficient individually.

## 7.2 Method of Steepest Descent

Consider the FIR Wiener filter of length  $M$ . We want to compute the filter coefficients iteratively.

Let us denote the time-varying filter parameters by

$$h_i[n], i = 0, 1, \dots, M-1$$

and define the filter parameter vector by

$$\mathbf{h}[n] = \begin{bmatrix} h_0[n] \\ h_1[n] \\ \vdots \\ h_{M-1}[n] \end{bmatrix}$$

We want to find the filter coefficients so as to minimize the mean-square error  $Ee^2[n]$

where

$$\begin{aligned} e[n] &= x[n] - \hat{x}[n] \\ &= x[n] - \sum_{i=0}^{M-1} h_i[n] y[n-i] \\ &= x[n] - \mathbf{h}'[n] \mathbf{y}[n] \\ &= x[n] - \mathbf{y}'[n] \mathbf{h}[n] \end{aligned}$$

$$\text{where } \mathbf{y}[n] = \begin{bmatrix} y[n] \\ y[n-1] \\ \vdots \\ y[n-M+1] \end{bmatrix}$$

Therefore

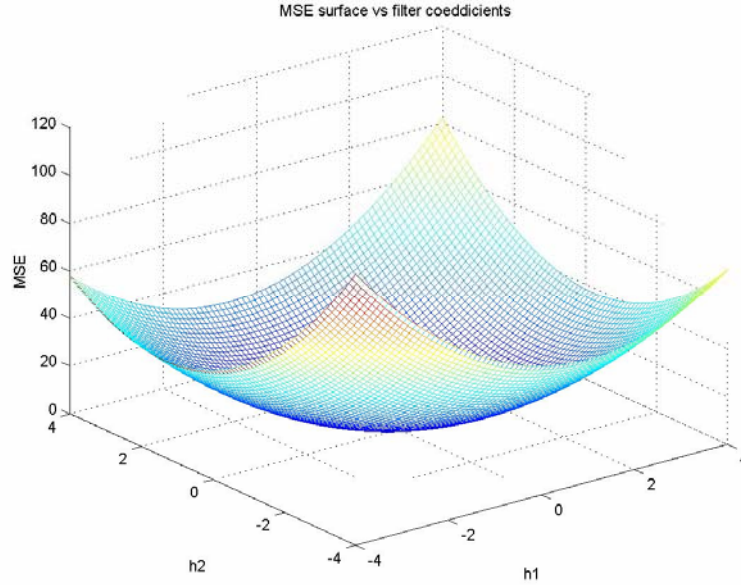
$$\begin{aligned} Ee^2[n] &= E(x[n] - \mathbf{h}'[n] \mathbf{y}[n])^2 \\ &= R_{xx}[0] - 2\mathbf{h}'[n] \mathbf{r}_{xy} + \mathbf{h}'[n] \mathbf{R}_{yy} \mathbf{h}[n] \end{aligned}$$

$$\text{where } \mathbf{r}_{xy} = \begin{bmatrix} R_{xy}[0] \\ R_{xy}[1] \\ \vdots \\ R_{xy}[M-1] \end{bmatrix}$$

and

$$\mathbf{R}_{yy} = \begin{bmatrix} R_{yy}[0] & R_{yy}[-1] & \dots & R_{yy}[1-M] \\ R_{yy}[1] & R_{yy}[0] & \dots & R_{yy}[2-M] \\ \dots & \dots & \dots & \dots \\ R_{yy}[M-1] & R_{yy}[M-2] & \dots & R_{yy}[0] \end{bmatrix}$$

- The cost function represented by  $Ee^2[n]$  is a quadratic in  $\mathbf{h}[n]$
- A unique global minimum exists
- The minimum is obtained by setting the gradient of  $Ee^2[n]$  to zero.



**Figure - Cost Function  $Ee^2[n]$**

The optimal set of filter parameters are given by

$$\mathbf{h}_{\text{opt}} = \mathbf{R}_{\text{XY}}^{-1} \mathbf{r}_{\text{XY}}$$

which is the FIR Wiener filter.

Many of the adaptive filter algorithms are obtained by simple modifications of the algorithms for deterministic optimization. Most of the popular adaptation algorithms are based on gradient-based optimization techniques, particularly the steepest descent technique.

The optimal Wiener filter can be obtained iteratively by the method of steepest descent. The optimum is found by updating the filter parameters by the rule

$$\mathbf{h}[n+1] = \mathbf{h}[n] + \frac{\mu}{2} (-\nabla Ee^2[n])$$

where

$$\begin{aligned} \nabla Ee^2[n] &= \begin{bmatrix} \frac{\partial Ee^2[n]}{\partial h_0} \\ \dots \\ \dots \\ \dots \\ \frac{\partial Ee^2[n]}{\partial h_{M-1}} \end{bmatrix} \\ &= -2\mathbf{r}_{\text{XY}} + 2\mathbf{R}_{\text{YY}}\mathbf{h}[n] \end{aligned}$$

and  $\mu$  is the step-size parameter.

So the steepest descent rule will now give

$$\mathbf{h}[n+1] = \mathbf{h}[n] + \mu(\mathbf{r}_{\text{XY}} - \mathbf{R}_{\text{YY}}\mathbf{h}[n])$$

### 7.3 Convergence of the steepest descent method

We have

$$\begin{aligned}\mathbf{h}[n+1] &= \mathbf{h}[n] + \mu(\mathbf{r}_{xy} - \mathbf{R}_{yy}\mathbf{h}[n]) \\ &= \mathbf{h}[n] - \mu\mathbf{R}_{yy}\mathbf{h}[n] + \mu\mathbf{r}_{xy} \\ &= (\mathbf{I} - \mu\mathbf{R}_{yy})\mathbf{h}[n] + \mu\mathbf{r}_{xy}\end{aligned}$$

where  $\mathbf{I}$  is the  $M \times M$  identity matrix.

This is a coupled set of linear difference equations.

Can we break it into simpler equations?

$\mathbf{R}_{yy}$  can be digitalized (KL transform) by the following relation

$$\mathbf{R}_{yy} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}'$$

where  $\mathbf{Q}$  is the orthogonal matrix of the eigenvectors of  $\mathbf{R}_{yy}$ .

$\mathbf{\Lambda}$  is a diagonal matrix with the corresponding eigen values as the diagonal elements.

$$\text{Also } \mathbf{I} = \mathbf{Q}\mathbf{Q}' = \mathbf{Q}'\mathbf{Q}$$

Therefore

$$\mathbf{h}[n+1] = (\mathbf{Q}\mathbf{Q}' - \mu\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}')\mathbf{h}[n] + \mu\mathbf{r}_{xy}$$

Multiply by  $\mathbf{Q}'$

$$\mathbf{Q}'\mathbf{h}[n+1] = (\mathbf{I} - \mu\mathbf{\Lambda})\mathbf{Q}'\mathbf{h}[n] + \mu\mathbf{Q}'\mathbf{r}_{xy}$$

**Define a new variable**

$$\bar{\mathbf{h}}[n] = \mathbf{Q}'\mathbf{h}[n] \text{ and } \bar{\mathbf{r}}_{xy} = \mathbf{Q}'\mathbf{r}_{xy}$$

Then

$$\begin{aligned}\bar{\mathbf{h}}[n+1] &= (\mathbf{I} - \mu\mathbf{\Lambda})\bar{\mathbf{h}}[n] + \mu\bar{\mathbf{r}}_{xy} \\ &= \begin{bmatrix} 1 - \mu\lambda_1 & 0 & \dots & \dots & 0 \\ 0 & & & & \\ \vdots & & & & \\ \vdots & & & & \\ 0 & \dots & \dots & & 1 - \mu\lambda_M \end{bmatrix} \bar{\mathbf{h}}[n] + \mu\bar{\mathbf{r}}_{xy}\end{aligned}$$

This is a decoupled set of linear difference equations

$$\bar{h}_i[n+1] = (1 - \mu\lambda_i)\bar{h}_i[n] + \mu\bar{r}_{xy}[i] \quad i = 0, 1, \dots, M-1$$

and can be easily solved for stability. The stability condition is given by

$$\begin{aligned}
& |1 - \mu\lambda_i| < 1 \\
\Rightarrow & -1 < 1 - \mu\lambda_i < 1 \\
\Rightarrow & 0 < \mu < 2 / \lambda_i, i = 1, \dots, M
\end{aligned}$$

Note that all the eigen values of  $\mathbf{R}_{YY}$  are positive.

Let  $\lambda_{\max}$  be the maximum eigen value. Then,

$$\begin{aligned}
\lambda_{\max} & < \lambda_1 + \lambda_2 + \dots + \lambda_M \\
& = \text{Trace}(\mathbf{R}_{YY}) \\
\therefore 0 < \mu & < \frac{2}{\text{Trace}(\mathbf{R}_{yy})} \\
& = \frac{2}{M \cdot R_{YY}[0]}
\end{aligned}$$

The steepest decent algorithm converges to the corresponding Wiener filter

$$\lim_{n \rightarrow \infty} \mathbf{h}[n] = \mathbf{R}_{YY}^{-1} \mathbf{r}_{XY}$$

if the step size  $\mu$  is within the range of specified by the above relation.]

## 7.4 Rate of Convergence

The rate of convergence of the Steepest Descent Algorithm will depend on the factor  $(1 - \mu\lambda_i)$  in

$$\bar{h}_i[n+1] = (1 - \mu\lambda_i) \bar{h}_i[n] + \mu \bar{R}_{xy}[i] \quad i = 0, 1, \dots, M-1$$

Thus the rate of convergence depends on the statistics of data and is related to the eigen value spread for the autocorrelation matrix. This rate is expressed using the condition

number of  $\mathbf{R}_{YY}$ , defined as  $k = \frac{\lambda_{\max}}{\lambda_{\min}}$  where  $\lambda_{\max}$  and  $\lambda_{\min}$  are respectively the maximum

and the minimum eigen values of  $\mathbf{R}_{YY}$ . The fastest convergence of this system occurs when  $k = 1$ , corresponding to white noise.

## 7.5 LMS algorithm (Least – Mean –Square) algorithm

Consider the steepest descent relation

$$\mathbf{h}[n+1] = \mathbf{h}[n] - \frac{\mu}{2} \nabla E e^2[n]$$

Where



$$\nabla Ee^2[n] = \begin{bmatrix} \frac{\partial Ee^2[n]}{\partial h_0} \\ \dots\dots\dots \\ \frac{\partial Ee^2[n]}{\partial h_{M-1}} \end{bmatrix}$$

In the LMS algorithm  $Ee^2[n]$  is approximated by  $e^2[n]$  to achieve a computationally simple algorithm.

$$\nabla Ee^2[n] \cong 2.e[n]. \begin{bmatrix} \frac{\partial e[n]}{\partial h_0} \\ \dots\dots\dots \\ \frac{\partial e[n]}{\partial h_{M-1}} \end{bmatrix}$$

Now consider

$$e[n] = x[n] - \sum_{i=0}^{M-1} h_i[n]y[n-i]$$

$$\frac{\partial e[n]}{\partial h_j} = -y[n-j], j = 0,1,\dots,M-1$$

$$\therefore \begin{bmatrix} \frac{\partial e[n]}{\partial h_0} \\ \dots\dots\dots \\ \frac{\partial e[n]}{\partial h_{M-1}} \end{bmatrix} = - \begin{bmatrix} y[n] \\ y[n-1] \\ \dots\dots\dots \\ y[n-M+1] \end{bmatrix} = -\mathbf{y}[n]$$

$$\therefore \nabla Ee^2[n] = -2e[n]\mathbf{y}[n]$$

The steepest descent update now becomes

$$\mathbf{h}[\mathbf{n} + 1] = \mathbf{h}[\mathbf{n}] + \mu e[n]\mathbf{y}[\mathbf{n}]$$

This modification is due to Widrow and Hopf and the corresponding adaptive filter is known as the **LMS filter**.

Hence the LMS algorithm is as follows

Given the input signal  $y[n]$ , reference signal  $x[n]$  and step size  $\mu$

1. Initialization  $h_i[0] = 0, i = 0,1,2,\dots,M-1$
- 2 For  $n > 0$

Filter output

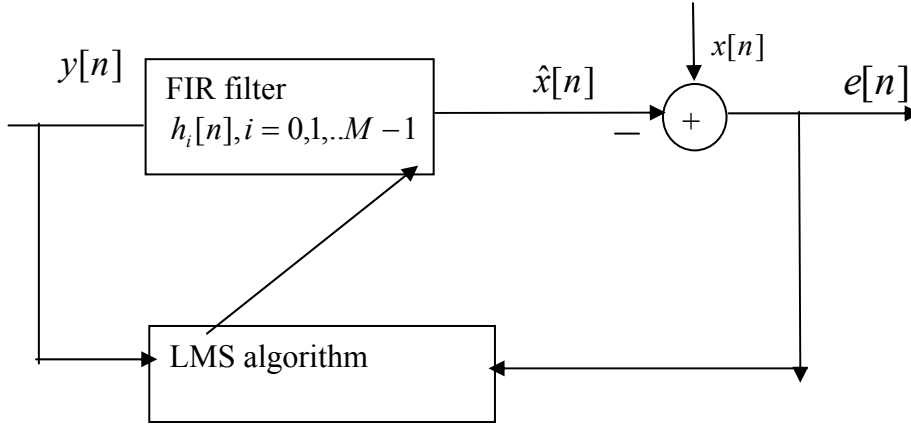
$$\hat{x}[n] = \mathbf{h}'[n]\mathbf{y}[n]$$

Estimation of the error

$$e[n] = x[n] - \hat{x}[n]$$

### 3. Tap weight adaptation

$$\mathbf{h}[\mathbf{n} + 1] = \mathbf{h}[\mathbf{n}] + \mu e[\mathbf{n}] \mathbf{y}[\mathbf{n}]$$



## 7.6 Convergence of the LMS algorithm

As there is a feedback loop in the adaptive algorithm, convergence is generally not assured. The convergence of the algorithm depends on the step size parameter  $\mu$ .

- The LMS algorithm is convergent in the mean if the step size parameter  $\mu$  satisfies the condition.

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

### Proof:

$$\mathbf{h}[\mathbf{n} + 1] = \mathbf{h}[\mathbf{n}] + \mu e[\mathbf{n}] \mathbf{y}[\mathbf{n}]$$

$$\begin{aligned} \therefore E\mathbf{h}[\mathbf{n} + 1] &= E\mathbf{h}[\mathbf{n}] + \mu E\mathbf{y}[\mathbf{n}]e[\mathbf{n}] \\ &= E\mathbf{h}[\mathbf{n}] + \mu E\mathbf{y}[\mathbf{n}](x[\mathbf{n}] - \mathbf{y}'[\mathbf{n}]\mathbf{h}[\mathbf{n}]) \\ &= E\mathbf{h}[\mathbf{n}] + \mu \mathbf{r}_{\mathbf{x}\mathbf{y}} - \mu E\mathbf{y}[\mathbf{n}]\mathbf{y}'[\mathbf{n}]\mathbf{h}[\mathbf{n}] \end{aligned}$$

Assuming the coefficient to be independent of data (Independence Assumption), we get

$$\begin{aligned} E\mathbf{h}[\mathbf{n} + 1] &= E\mathbf{h}[\mathbf{n}] + \mu \mathbf{r}_{\mathbf{x}\mathbf{y}} - \mu E\mathbf{y}[\mathbf{n}]\mathbf{y}'[\mathbf{n}]E\mathbf{h}[\mathbf{n}] \\ &= E\mathbf{h}[\mathbf{n}] + \mu \mathbf{r}_{\mathbf{x}\mathbf{y}} - \mu \mathbf{R}_{\mathbf{y}\mathbf{y}}E\mathbf{h}[\mathbf{n}] \end{aligned}$$

Hence the mean value of the filter coefficients satisfies the steepest descent iterative relation so that the same stability condition applies to the mean of the filter coefficients.

- In the practical situation, knowledge of  $\lambda_{\max}$  is not available and Trace  $\mathbf{R}_{\mathbf{y}\mathbf{y}}$  can be taken as the conservative estimate of  $\lambda_{\max}$  so that for convergence
- $0 < \mu < \frac{2}{\text{Trace}(\mathbf{R}_{\mathbf{y}\mathbf{y}})}$
- Also note that trace,  $\text{Trace}(\mathbf{R}_{\mathbf{y}\mathbf{y}}) = M\mathbf{R}_{\mathbf{y}\mathbf{y}}[0] = \text{Tape input power of the LMS filter.}$

Generally, a too small value of  $\mu$  results in slower convergence where as big values of  $\mu$  will result in larger fluctuations from the mean. Choosing a proper value of  $\mu$  is very important for the performance of the LMS algorithm.

In addition, the rate of convergence depends on the statistics of data and is related to the eigenvalue spread for the autocorrelation matrix. This is defined using the condition number of  $\mathbf{R}_{YY}$ , defined as  $k = \frac{\lambda_{\max}}{\lambda_{\min}}$  where  $\lambda_{\min}$  is the minimum eigenvalue of  $\mathbf{R}_{YY}$ . The

fastest convergence of this system occurs when  $k = 1$ , corresponding to white noise. This states that the fastest way to train a LMS adaptive system is to use white noise as the training input. As the noise becomes more and more colored, the speed of the training will decrease.

The average of each filter tap –weight converges to the corresponding optimal filter tap-weight. But this does not ensure that the coefficients converge to the optimal values.

## 7.7 Excess mean square error

Consider the LMS difference equation:

$$\mathbf{h}[n+1] = \mathbf{h}[n] + \mu e[n] \mathbf{y}[n]$$

We have seen that the mean of LMS coefficient converges to the steepest descent solution. But this does not guarantee that the mean square error of the LMS estimator will converge to the mean square error corresponding to the wiener solution. There is a fluctuation of the LMS coefficient from the wiener filter coefficient.

Let  $\mathbf{h}_{\text{opt}}$  = optimal wiener filter impulse response.

The instantaneous deviation of the LMS coefficient from  $\mathbf{h}_{\text{opt}}$  is

$$\Delta \mathbf{h}[n] = \mathbf{h}[n] - \mathbf{h}_{\text{opt}}$$

$$\begin{aligned} \varepsilon[n] &= E e^2[n] = E \{x[n] - \mathbf{h}'_{\text{opt}} \mathbf{y}[n] - \Delta \mathbf{h}'[n] \mathbf{y}[n]\}^2 \\ &= E \{x[n] - \mathbf{h}'_{\text{opt}} \mathbf{y}[n]\}^2 + E \Delta \mathbf{h}'[n] \mathbf{y}[n] \mathbf{y}'[n] \Delta \mathbf{h}[n] - 2E(e_{\text{opt}}[n] \Delta \mathbf{h}'[n] \mathbf{y}[n]) \\ &= \varepsilon_{\min} + E \Delta \mathbf{h}'[n] \mathbf{y}[n] \mathbf{y}'[n] \Delta \mathbf{h}[n] - 2E(e_{\text{opt}}[n] \Delta \mathbf{h}'[n] \mathbf{y}[n]) \\ &= \varepsilon_{\min} + E \Delta \mathbf{h}'[n] \mathbf{y}[n] \mathbf{y}'[n] \Delta \mathbf{h}[n] \end{aligned}$$

assuming the independence of deviation with respect to data and at  $E \Delta \mathbf{h}[n] = 0$ .

Therefore,

$$\varepsilon_{\text{excess}} = E \Delta \mathbf{h}'[n] \mathbf{y}[n] \mathbf{y}'[n] \Delta \mathbf{h}[n]$$

An exact analysis of the excess mean-square error is quite complicated and its approximate value is given by

$$\varepsilon_{\text{excess}} = \varepsilon_{\min} \frac{\sum_{i=1}^M \frac{\mu \lambda_i}{2 - \mu \lambda_i}}{1 - \sum_{i=1}^M \frac{\mu \lambda_i}{2 - \mu \lambda_i}}$$

The LMS algorithm is said to converge in the mean-square sense provided the step-length parameter satisfies the relations

$$\mu \sum_{i=1}^M \frac{\mu \lambda_i}{2 - \mu \lambda_i} < 1$$

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

If  $\sum_{i=1}^M \frac{\mu \lambda_i}{2 - \mu \lambda_i} \ll 1$

then  $\varepsilon_{\text{excess}} = \varepsilon_{\min} \sum_{i=1}^M \frac{\mu \lambda_i}{2 - \mu \lambda_i}$

Further, if

$$\begin{aligned} \mu &\ll 1 \\ \varepsilon_{\text{excess}} &= \varepsilon_{\min} \mu \frac{\frac{1}{2} \text{Trace}(\mathbf{R}_{\mathbf{Y}\mathbf{Y}})}{1 - 0} \\ &\simeq \varepsilon_{\min} \mu \frac{1}{2} \text{Trace}(\mathbf{R}_{\mathbf{Y}\mathbf{Y}}) \end{aligned}$$

The factor  $\frac{\varepsilon_{\text{excess}}}{\varepsilon_{\min}} = \sum_{i=1}^M \frac{\mu \lambda_i}{2 - \mu \lambda_i}$  is called the *misadjustment factor* for the LMS filter.

## 7.8 Drawback of the LMS Algorithm

Convergence is slow when the eigenvalue spread of the autocorrelation matrix is large.

Misadjustment factor given by

$$\frac{\varepsilon_{\text{excess}}}{\varepsilon_{\min}} \approx \frac{1}{2} \mu \text{Trace}(\mathbf{R}_{\mathbf{Y}\mathbf{Y}})$$

is large unless  $\mu$  is much smaller. Thus the selection of the step-size parameter is crucial in the case of the LMS algorithm.

When the input signal is nonstationary the eigenvalues also change with time and selection of  $\mu$  becomes more difficult.

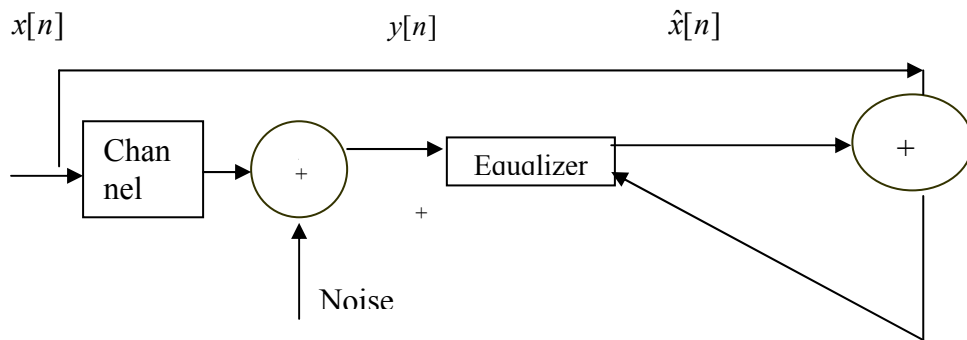
**Example 1:**

The input to a communication channel is a test sequence  $x[n] = 0.8x[n-1] + w[n]$  where  $w[n]$  is a 0 mean unity variance white noise. The channel transfer function is given by  $H(z) = z^{-1} - 0.5z^{-2}$  and the channel is affected white Gaussian noise of variance 1.

- Find the FIR Wiener filter of length 2 for channel equalization
- Write down the LMS filter update equations
- Find the bounds of LMS step length parameters
- Find the excess mean square error.

**Solution:**

From the given model



$$R_{xx}[m] = \frac{\sigma_w^2}{1 - 0.8^2} (0.8)^{|m|}$$

$$R_{xx}[0] = 2.78$$

$$R_{xx}[1] = 2.22$$

$$R_{xx}[2] = 1.78$$

$$R_{xx}[3] = 1.42$$

$$y[n] = x[n-1] - 0.5x[n-2] + v(n)$$

$$\therefore R_{yy}[m] = 1.25R_{xx}[m] - 0.5R_{xx}[m+1] - 0.5R_{xx}[m-1] + \delta[m]$$

$$R_{yy}[0] = 2.255$$

$$R_{yy}[1] = 0.5325$$

$$\text{also } R_{xy}[m] = Ex[n](x[n-1-m] - 0.5x[n-2-m] + v[n-m])$$

$$= R_{xx}[m+1] - 0.5R_{xx}[m+2]$$

$$\therefore R_{xy}[0] = 1.33$$

$$\text{and } R_{xy}[1] = 1.07$$

Therefore, the Wiener solution is given by

$$\begin{bmatrix} 2.255 & 0.533 \\ 0.533 & 2.255 \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \end{bmatrix} = \begin{bmatrix} 1.33 \\ 1.07 \end{bmatrix}$$

$$MSE = R_{xx}[0] - h_0 R_{xy}[1] - h_1 R_{xy}[2]$$

$$h_0 = 0.51$$

$$h_1 = 0.35$$

$$\lambda_1, \lambda_2 = 2.79, 1.72$$

$$\mu < \frac{2}{2.79}$$

$$= 0.72$$

Excess mean square error

$$= \xi_{mn} \frac{\sum_{i=1}^2 \frac{2\lambda_i}{2 - \mu\lambda_i}}{1 - \sum_{i=1}^2 \frac{\mu\lambda_i}{2 - 2\lambda_i}}$$

## 7.9 Leaky LMS Algorithm

Minimizes  $e^2[n] + \alpha \|\mathbf{h}[n]\|^2$

where  $\|\mathbf{h}[n]\|$  is the modulus of the LMS weight vector and  $\alpha$  is a positive quantity.

The corresponding algorithm is given by

$$\mathbf{h}[n+1] = (1 - \mu\alpha)\mathbf{h}[n] + \mu\alpha\mathbf{e}[n]\mathbf{y}[n]$$

where  $\mu\alpha$  is chosen to be less than 1. In such a situation the pole will be inside the unit circle, instability problem will not be there and the algorithm will converge.

## 7.10 Normalized LMS Algorithm

For convergence of the LMS algorithm

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

and the conservative bound is given by

$$0 < \mu < \frac{2}{\text{Trace}(\mathbf{R}_{YY})}$$

$$= \frac{2}{M\mathbf{R}_{YY}[0]}$$

$$= \frac{2}{ME(Y^2[n])}$$

We can estimate the bound by estimating  $E(Y^2[n])$  by

$\frac{1}{M} \sum_{n=0}^M Y^2[n]$  so that we get the bound

$$0 < \mu < \frac{2}{\sum_{i=0}^{M-1} y^2[n-i]} = \frac{2}{\|\mathbf{y}[n]\|^2}$$

Then we can take

$$\mu = \beta \frac{2}{\|\mathbf{y}[n]\|^2}$$

and the LMS updating becomes

$$\mathbf{h}[n+1] = \mathbf{h}[n] + \beta \frac{1}{\|\mathbf{y}[n]\|^2} e[n] \mathbf{y}[n]$$

where  $0 < \beta < 2$

## 7.11 Discussion - LMS

The NLMS algorithm has more computational complexity compared to the LMS algorithm

Under certain assumptions, it can be shown that the NLMS algorithm converges for  $0 < \beta < 2$

$\|\mathbf{y}[n]\|^2$  can be efficiently estimated using the recursive relation

$$\|\mathbf{y}[n]\|^2 = \|\mathbf{y}[n-1]\|^2 + y^2[n] - y^2[n-M-1]$$

Notice that the NLMS algorithm does not change the direction of updation in steepest descent algorithm

If  $\mathbf{y}[n]$  is close to zero, the denominator term ( $\|\mathbf{y}[n]\|^2$ ) in NLMS equation becomes very small and

$$\mathbf{h}[n+1] = \mathbf{h}[n] + \beta \frac{1}{\|\mathbf{y}[n]\|^2} e[n] \mathbf{y}[n] \quad \text{may diverge}$$

To overcome this drawback a small positive number  $\varepsilon$  is added to the denominator term the NLMS equation. Thus

$$\mathbf{h}[n+1] = \mathbf{h}[n] + \beta \frac{1}{\varepsilon + \|\mathbf{y}[n]\|^2} e[n] \mathbf{y}[n]$$

For computational efficiency, other modifications are suggested to the LMS algorithm. Some of the modified algorithms are *blocked-LMS* algorithm, *signed LMS* algorithm etc. LMS algorithm can be obtained for IIR filter to adaptively update the parameters of the filter

$$y[n] = \sum_{i=1}^{M-1} a_i[n] y[n-i] + \sum_{i=0}^{N-1} b_i[n] x[n-i]$$

How ever, IIR LMS algorithm has poor performance compared to FIR LMS filter.



## 7.12 Recursive Least Squares (RLS) Adaptive Filter

- LMS convergence slow
- Step size parameter is to be properly chosen
- Excess mean-square error is high
- LMS minimizes the instantaneous square error  $e^2[n]$
- Where  $e[n] = x[n] - \mathbf{h}'[n] \mathbf{y}[n] = x[n] - \mathbf{y}'[n] \mathbf{h}[n]$

The RLS algorithm considers all the available data for determining the filter parameters.

The filter should be optimum with respect to all the available data in certain sense.

Minimizes the cost function

$$\varepsilon[n] = \sum_{k=0}^n \lambda^{n-k} e^2[k]$$

with respect to the filter parameter vector  $\mathbf{h}[n] = \begin{bmatrix} h_0[n] \\ h_1[n] \\ \vdots \\ h_{M-1}[n] \end{bmatrix}$

where  $\lambda$  is the weighing factor known as the forgetting factor

- Recent data is given more weightage
- For stationary case  $\lambda = 1$  can be taken
- $\lambda \cong 0.99$  is effective in tracking local nonstationarity

The minimization problem is

Minimize  $\varepsilon[n] = \sum_{k=0}^n \lambda^{n-k} (x[k] - \mathbf{y}'[k] \mathbf{h}[n])^2$  with respect to  $\mathbf{h}[n]$

The minimum is given by

$$\frac{\partial \varepsilon(n)}{\partial \mathbf{h}(n)} = \mathbf{0}$$

$$\Rightarrow 2 \sum_{k=0}^n \lambda^{n-k} (x[k] \mathbf{y}[k] - \mathbf{y}[k] \mathbf{y}'[k] \mathbf{h}[n]) = 0$$

$$\Rightarrow \mathbf{h}[n] = \left( \sum_{k=0}^n \lambda^{n-k} \mathbf{y}[k] \mathbf{y}'[k] \right)^{-1} \sum_{k=0}^n \lambda^{n-k} x[k] \mathbf{y}[k]$$

Let us define  $\hat{\mathbf{R}}_{YY}[n] = \sum_{k=0}^n \lambda^{n-k} \mathbf{y}[k] \mathbf{y}'[k]$

which is an estimator for the autocorrelation matrix  $\mathbf{R}_{YY}$ .

Similarly  $\hat{\mathbf{r}}_{XY}[n] = \sum_{k=0}^n \lambda^{n-k} x[k] \mathbf{y}[k] =$  estimator for the autocorrelation vector  $\mathbf{r}_{XY}[n]$

Hence  $\mathbf{h}[n] = (\hat{\mathbf{R}}_{XY}[n])^{-1} \hat{\mathbf{r}}_{XY}[n]$

Matrix inversion is involved which makes the direct solution difficult. We look forward for a recursive solution.

### 7.13 Recursive representation of $\hat{\mathbf{R}}_{YY}[n]$

$\hat{\mathbf{R}}_{YY}[n]$  can be rewritten as follows

$$\begin{aligned} \hat{\mathbf{R}}_{YY}[n] &= \sum_{k=0}^{n-1} \lambda^{n-k} \mathbf{y}[k] \mathbf{y}'[k] + \mathbf{y}[n] \mathbf{y}'[n] \\ &= \lambda \sum_{k=0}^{n-1} \lambda^{n-1-k} \mathbf{y}[k] \mathbf{y}'[k] + \mathbf{y}[n] \mathbf{y}'[n] \\ &= \lambda \hat{\mathbf{R}}_{YY}[n-1] + \mathbf{y}[n] \mathbf{y}'[n] \end{aligned}$$

This shows that the autocorrelation matrix can be recursively computed from its previous values and the present data vector.

Similarly  $\hat{\mathbf{r}}_{XY}[n] = \lambda \hat{\mathbf{r}}_{XY}[n-1] + x[n] \mathbf{y}[n]$

$$\begin{aligned} \mathbf{h}[n] &= [\hat{\mathbf{R}}_{YY}[n]]^{-1} \hat{\mathbf{r}}_{XY}[n] \\ &= (\lambda \hat{\mathbf{R}}_{YY}[n-1] + \mathbf{y}[n] \mathbf{y}'[n])^{-1} \hat{\mathbf{r}}_{XY}[n] \end{aligned}$$

For the matrix inversion above the matrix inversion lemma will be useful.

### 7.14 Matrix Inversion Lemma

If  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$  and  $\mathbf{D}$  are matrices of proper orders,  $\mathbf{A}$  and  $\mathbf{C}$  nonsingular

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{D} \mathbf{A}^{-1} \mathbf{B} + \mathbf{C}^{-1})^{-1} \mathbf{D} \mathbf{A}^{-1}$$

Taking  $\mathbf{A} = \lambda \hat{\mathbf{R}}_{YY}[n-1]$ ,  $\mathbf{B} = \mathbf{y}[n]$ ,  $\mathbf{C} = 1$  and  $\mathbf{D} = \mathbf{y}'[n]$

we will have

$$\begin{aligned} (\hat{\mathbf{R}}_{YY}[n])^{-1} &= \frac{1}{\lambda} \hat{\mathbf{R}}_{YY}^{-1}[n-1] - \frac{1}{\lambda} \hat{\mathbf{R}}_{YY}^{-1}[n-1] \mathbf{y}[n] \left( \mathbf{y}'[n] \frac{1}{\lambda} [\hat{\mathbf{R}}_{YY}^{-1}[n-1] \mathbf{y}[n] + 1] \right)^{-1} \mathbf{y}'[n] \hat{\mathbf{R}}_{YY}^{-1}[n-1] \\ &= \frac{1}{\lambda} \left( \hat{\mathbf{R}}_{YY}^{-1}[n-1] - \frac{\hat{\mathbf{R}}_{YY}^{-1}[n-1] \mathbf{y}[n] \mathbf{y}'[n] \hat{\mathbf{R}}_{YY}^{-1}[n-1]}{\lambda + \mathbf{y}'[n] \hat{\mathbf{R}}_{YY}^{-1}[n-1] \mathbf{y}[n]} \right) \end{aligned}$$

Rename  $\mathbf{P}[n] = \hat{\mathbf{R}}_{YY}^{-1}[n]$ . Then

$$\mathbf{P}[n] = \frac{1}{\lambda} (\mathbf{P}[n-1] - \mathbf{k}[n] \mathbf{y}'[n] \mathbf{P}[n-1])$$

where  $\mathbf{k}[n]$  is called the ‘gain vector’ and given by

$$\mathbf{k}[n] = \frac{\mathbf{P}[n-1] \mathbf{y}[n]}{\lambda + \mathbf{y}'[n] \mathbf{P}[n-1] \mathbf{y}[n]}$$

$\mathbf{k}[n]$  important to interpret adaptation is also related to the current data vector  $\mathbf{y}[n]$  by

$$\mathbf{k}[n] = \mathbf{P}[n] \mathbf{y}[n]$$

To establish the above relation consider

$$\mathbf{P}[n] = \frac{1}{\lambda} (\mathbf{P}[n-1] - \mathbf{k}[n] \mathbf{y}'[n] \mathbf{P}[n-1])$$

Multiplying by  $\lambda$  and post-multiplying by  $\mathbf{y}[n]$  and simplifying we get

$$\begin{aligned} \lambda \mathbf{P}[n] \mathbf{y}[n] &= (\mathbf{P}[n-1] - \mathbf{k}[n] \mathbf{y}'[n] \mathbf{P}[n-1]) \mathbf{y}[n] \\ &= \mathbf{P}[n-1] \mathbf{y}[n] - \mathbf{k}[n] \mathbf{y}'[n] \mathbf{P}[n-1] \mathbf{y}[n] \\ &= \lambda \mathbf{k}[n] \end{aligned}$$

Therefore

$$\begin{aligned} \mathbf{h}[n] &= (\hat{\mathbf{R}}_{\mathbf{Y}\mathbf{Y}}[n])^{-1} \hat{\mathbf{r}}_{\mathbf{X}\mathbf{Y}}[n] \\ &= \mathbf{P}[n] (\lambda \hat{\mathbf{r}}_{\mathbf{X}\mathbf{Y}}[n-1] + x[n] \mathbf{y}[n]) \\ &= \lambda \mathbf{P}[n] \hat{\mathbf{r}}_{\mathbf{X}\mathbf{Y}}[n-1] + x[n] \mathbf{P}[n] \mathbf{y}[n] \\ &= \lambda \frac{1}{\lambda} [\mathbf{P}[n-1] - \mathbf{k}[n] \mathbf{y}'[n] \mathbf{P}[n-1]] \hat{\mathbf{r}}_{\mathbf{X}\mathbf{Y}}[n-1] + x[n] \mathbf{P}[n] \mathbf{y}[n] \\ &= \mathbf{h}[n-1] - \mathbf{k}[n] \mathbf{y}'[n] \mathbf{h}[n-1] + x[n] \mathbf{k}[n] \\ &= \mathbf{h}[n-1] + \mathbf{k}[n] (x[n] - \mathbf{y}'[n] \mathbf{h}[n-1]) \end{aligned}$$

## 7.15 RLS algorithm Steps

### Initialization:

At  $n = 0$

$$\begin{aligned} \mathbf{P}[0] &= \delta \mathbf{I}_{\mathbf{M}\times\mathbf{M}}, \quad \delta \text{ a postive number} \\ \mathbf{y}[0] &= \mathbf{0}, \quad \mathbf{h}[0] = \mathbf{0} \end{aligned}$$

Choose  $\lambda$

### Operation:

For 1 to  $n = \text{Final}$  do

1. Get  $x[n], \mathbf{y}[n]$
2. Get  $e[n] = x[n] - \mathbf{h}'[n-1] \mathbf{y}[n]$
3. Calculate gain vector  $\mathbf{k}[n] = \frac{\mathbf{P}[n-1] \mathbf{y}[n]}{\lambda + \mathbf{y}'[n] \mathbf{P}[n-1] \mathbf{y}[n]}$

4. Update the filter parameters

$$\mathbf{h}[n] = \mathbf{h}[n-1] + \mathbf{k}[n]e[n]$$

5. Update the  $\mathbf{P}$  matrix

$$\mathbf{P}[n] = \frac{1}{\lambda} (\mathbf{P}[n-1] - \mathbf{k}[n]\mathbf{y}'[n]\mathbf{P}[n-1])$$

end do

## 7.16 Discussion – RLS

### 7.16.1 Relation with Wiener filter

We have the optimality condition analysis for the RLS filters

$$\hat{\mathbf{R}}_{\mathbf{Y}\mathbf{Y}}[n]\mathbf{h}[n] = \hat{\mathbf{r}}_{\mathbf{Y}\mathbf{X}}[n]$$

where  $\hat{\mathbf{R}}_{\mathbf{Y}\mathbf{Y}}[n] = \sum_{k=0}^n \lambda^{n-k} \mathbf{y}[k]\mathbf{y}'[k]$

Dividing by  $n+1$

$$\frac{\hat{\mathbf{R}}_{\mathbf{Y}\mathbf{Y}}[n]}{n+1} = \frac{\sum_{k=0}^n \lambda^{n-k} \mathbf{y}[k]\mathbf{y}'[k]}{n+1}$$

if we consider the elements of  $\frac{\hat{\mathbf{R}}_{\mathbf{Y}\mathbf{Y}}[n]}{n+1}$ , we see that each is an estimator for the auto-correlation of specific lag.

$$\lim_{n \rightarrow \infty} \frac{\hat{\mathbf{R}}_{\mathbf{Y}\mathbf{Y}}[n]}{n+1} = \mathbf{R}_{\mathbf{Y}\mathbf{Y}}[n]$$

in the mean square sense. In other words, weighted sample autocorrelation is a consistent estimator of the auto-correlation function of a WSS process.

Similarly

$$\lim_{n \rightarrow \infty} \frac{\hat{\mathbf{r}}_{\mathbf{Y}\mathbf{X}}[n]}{n+1} = \mathbf{r}_{\mathbf{Y}\mathbf{X}}[n]$$

Hence as  $n \rightarrow \infty$ , optimality condition can be written as

$$\mathbf{R}_{\mathbf{Y}\mathbf{Y}}[n]\mathbf{h}[n] = \mathbf{r}_{\mathbf{Y}\mathbf{X}}[n].$$

### 7.16.2. Dependence condition on the initial values

Consider the recursive relation

$$\hat{\mathbf{R}}_{YY}[n] = \lambda \hat{\mathbf{R}}_{YY}[n-1] + \mathbf{y}[n]\mathbf{y}'[n]$$

Corresponding to

$$\hat{\mathbf{R}}_{YY}^{-1}[0] = \delta \mathbf{I}$$

we have  $\hat{\mathbf{R}}_{YY}[0] = \frac{\mathbf{I}}{\delta}$

With this initial condition the matrix difference equation has the solution

$$\begin{aligned} \tilde{\mathbf{R}}[n] &= \lambda^{n+1} \hat{\mathbf{R}}_{YY}[-1] + \sum_{k=0}^n \lambda^{n-k} \mathbf{y}[k]\mathbf{y}'[k] \\ &= \lambda^{n+1} \hat{\mathbf{R}}_{YY}[-1] + \hat{\mathbf{R}}_{YY}[n] \\ &= \lambda^{n+1} \frac{\mathbf{I}}{\delta} + \hat{\mathbf{R}}_{YY}[n] \end{aligned}$$

Hence the optimality condition is modified as

$$(\lambda^{n+1} \frac{\mathbf{I}}{\delta} + \hat{\mathbf{R}}_{YY}[n]) \tilde{\mathbf{h}}[n] = \hat{\mathbf{r}}_{XY}[n]$$

where  $\tilde{\mathbf{h}}[n]$  is the modified solution due to assumed initial value of the P-matrix.

$$\frac{\lambda^{n+1} \hat{\mathbf{R}}_{YY}^{-1}[n] \tilde{\mathbf{h}}[n]}{\delta} + \tilde{\mathbf{h}}[n] = \mathbf{h}[n]$$

If we take  $\lambda$  as less than 1, then the bias term in the left-hand side of the above equation will be eventually die down and we will get

$$\tilde{\mathbf{h}}[n] = \mathbf{h}[n]$$

### 7.16.3. Convergence in stationary condition

- If the data is stationary, the algorithm will converge in mean at best in M iterations, where M is the number of taps in the adaptive FIR filter.
- The filter coefficients converge in the mean to the corresponding Wiener filter coefficients.
- Unlike the LMS filter which converges in the mean at infinite iterations, the RLS filter converges in a finite time. Convergence is less sensitive to eigen value spread. This is a remarkable feature of the RLS algorithm.
- The RLS filter can also be shown to converge to the Wiener filter in the mean-square sense so that there is zero excess mean-square error.

#### 7.16.4. Tracking non-stationarity

If  $\lambda$  is small  $\lambda^{n-i} \cong 0$  for  $i \ll n$

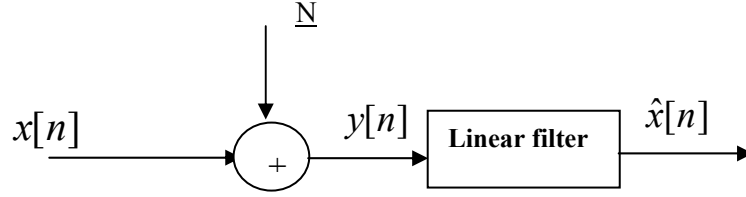
$\Rightarrow$  the filter is based on most recent values. This also qualitatively explains that the filter can track non stationary in data.

#### 7.16.5. Computational Complexity

Several matrix multiplications result in  $\cong 7M^2$  arithmetic operations, which are quite large, if the filter tap-length is high. So we have to go for the best implementation of the RLS algorithm.

## CHAPTER – 8: KALMAN FILTER

### 8.1 Introduction



To estimate a signal  $x[n]$  in the presence of noise,

- FIR Wiener Filter is optimum when the data length and the filter length are equal.
- IIR Wiener Filter is based on the assumption that infinite length of data sequence is available.

Neither of the above filters represents the physical situation. We need a filter that adds a tap with each addition of data.

The basic mechanism in Kalman filter is to estimate the signal recursively by the following relation

$$\hat{x}[n] = A_n \hat{x}[n-1] + K_n y[n]$$

The whole of Kalman filter is also based on the innovation representation of the signal. We used this model to develop causal IIR Wiener filter.

### 8.2 Signal Model

The simplest Kalman filter uses the first-order AR signal model

$$x[n] = ax[n-1] + w[n]$$

where  $w[n]$  is a white noise sequence.

The observed data is given by

$$y[n] = x[n] + v[n]$$

where  $v[n]$  is another white noise sequence independent of the signal.

The general stationary signal is modeled by a difference equation representing the ARMA (p,q) model. Such a signal can be modeled by the state-space model and is given by

$$\mathbf{x}[n] = \mathbf{A}\mathbf{x}[n-1] + \mathbf{B}w[n] \quad (1)$$

And the observations can be represented as a linear combination of the 'states' and the observation noise.

$$y[n] = \mathbf{c}'\mathbf{x}[n] + v[n] \quad (2)$$

Equations (1) and (2) have direct relation with the state space model in the control system where you have to estimate the ‘unobservable’ states of the system through an observer that performs well against noise.

**Example 1:**

Consider the  $AR(M)$  model

$$x[n] = a_1 x[n-1] + a_2 x[n-2] + \dots + a_M x[n-M] + w[n]$$

Then the state variable model for  $x[n]$  is given by

$$\mathbf{x}[n] = \mathbf{A}\mathbf{x}[n-1] + \mathbf{B}w[n]$$

where

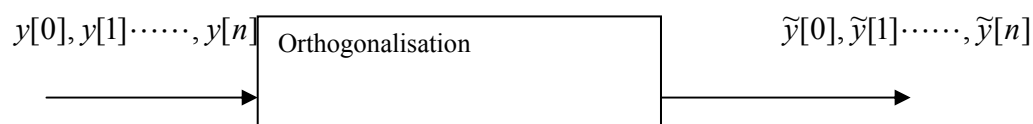
$$\mathbf{x}[n] = \begin{bmatrix} x_1[n] \\ x_2[n] \\ \vdots \\ x_M[n] \end{bmatrix}, \quad x_1[n] = x[n], \quad x_2[n] = x[n-1], \dots \text{ and } x_M[n] = x[n-M+1],$$

$$\mathbf{A} = \begin{bmatrix} a_1 & a_2 & \dots & a_M \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

$$\text{and } \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Our analysis will include only the simple (scalar) Kalman filter

The Kalman filter also uses the innovation representation of the stationary signal as does by the IIR Wiener filter. The innovation representation is shown in the following diagram.



Let  $\hat{x}[n]$  be the LMMSE of  $x[n]$  based on the data  $y[0], y[1], \dots, y[n]$ .



In the above representation  $\tilde{y}[n]$  is the innovation of  $y[n]$  and contains the same information as the original sequence. Let  $\hat{E}(y[n] | y[n-1], \dots, y[0])$  be the linear prediction of  $y[n]$  based on  $y[n-1], \dots, y[0]$ .

Then

$$\begin{aligned}
 \tilde{y}[n] &= y[n] - \hat{E}(y[n] | y[n-1], \dots, y[0]) \\
 &= y[n] - \hat{E}(x[n] + v[n] | y[n-1], \dots, y[0]) \\
 &= y[n] - \hat{E}(ax[n-1] + w[n] + v[n] | y[n-1], \dots, y[0]) \\
 &= y[n] - \hat{E}(ax[n-1] | y[n-1], \dots, y[0]) \\
 &= y[n] - a\hat{x}[n-1] \\
 &= y[n] - x[n] + x[n] - a\hat{x}[n-1] \\
 &= y[n] - x[n] + ax[n-1] + w[n] - a\hat{x}[n-1] \\
 &= y[n] - x[n] + a(x[n-1] - \hat{x}[n-1]) + w[n] \\
 &= v[n] + ae[n-1] + w[n]
 \end{aligned}$$

which is a linear combination of three mutually independent orthogonal sequences.

It can be easily shown that  $\tilde{y}[n]$  is orthogonal to each of  $\tilde{y}[n-1], \tilde{y}[n-2], \dots$  and  $\tilde{y}[0]$ .

The LMMSE estimation of  $x[n]$  based on  $y[0], y[1], \dots, y[n]$ , is same as the estimation based on the innovation sequence,  $\tilde{y}[0], \tilde{y}[1], \dots, \tilde{y}[n-1], \tilde{y}[n]$ . Therefore,

$$\hat{x}[n] = \sum_{i=0}^n k_i \tilde{y}[i]$$

where  $k_i$  can be obtained by the orthogonality relation.

Consider the relation

$$\begin{aligned}
 x[n] &= \hat{x}[n] + e[n] \\
 &= \sum_{i=0}^n k_i \tilde{y}[i] + e[n]
 \end{aligned}$$

Then

$$E(x[n] - \sum_{i=0}^n k_i \tilde{y}[i]) \tilde{y}[j] = 0 \quad j = 0, 1, \dots, n$$

so that

$$k_j = E x[n] \tilde{y}[j] / \sigma_j^2 \quad j = 0, 1, \dots, n$$

Similarly

$$\begin{aligned}
 x[n-1] &= \hat{x}[n-1] + e[n-1] \\
 &= \sum_{i=0}^{n-1} k'_i \tilde{y}[i] + e[n-1] \\
 k'_j &= Ex[n-1]\tilde{y}[j] / \sigma_j^2 \quad j = 0, 1..n-1 \\
 &= E(x[n] - w[n])\tilde{y}[j] / a\sigma_j^2 \quad j = 0, 1..n-1 \\
 &= E(x[n])\tilde{y}[j] / a\sigma_j^2 \quad j = 0, 1..n-1 \\
 \therefore k'_j &= k_j / a \quad j = 0, 1..n-1
 \end{aligned}$$

Again,

$$\begin{aligned}
 \hat{x}[n] &= \sum_{i=0}^n k_i \tilde{y}[i] = \sum_{i=0}^{n-1} k_i \tilde{y}[i] + k_n \tilde{y}[n] \\
 &= a \sum_{i=0}^{n-1} k'_i \tilde{y}[i] + k_n \tilde{y}[n] \\
 &= a\hat{x}[n-1] + k_n (y[n] - \hat{E}(y[n] / y[n-1], \dots, y[0])) \\
 &= a\hat{x}[n-1] + k_n (y[n] - a\hat{x}[n-1]) \\
 &= (1 - k_n)a\hat{x}[n-1] + k_n y[n]
 \end{aligned}$$

where  $\hat{E}(y[n] / y[n-1], \dots, y[0])$  is the linear prediction of  $y[n]$  based on observations  $y[n-1], \dots, y[0]$  and which is same as  $a\hat{x}[n-1]$ .

$$\begin{aligned}
 \therefore \hat{x}[n] &= A_n \hat{x}[n-1] + k_n y[n] \\
 \text{with } A_n &= (1 - k_n)a
 \end{aligned}$$

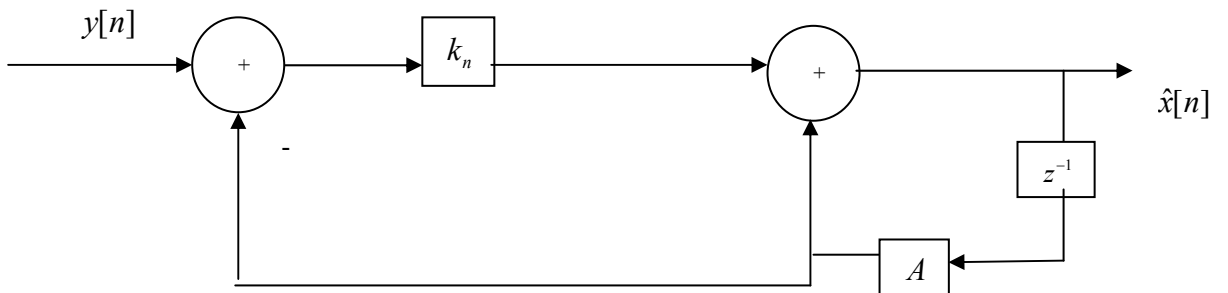
Thus the recursive estimator  $\hat{x}[n]$  is given by

$$\hat{x}[n] = A_n \hat{x}[n-1] + k_n y[n]$$

Or

$$\hat{x}[n] = a\hat{x}[n-1] + k_n (y[n] - a\hat{x}[n-1])$$

The filter can be represented in the following diagram



### 8.3 Estimation of the filter-parameters

Consider the estimator

$$\hat{x}[n] = A_n \hat{x}[n-1] + k_n y[n]$$

The estimation error is given by

$$e[n] = x[n] - \hat{x}[n]$$

Therefore  $e[n]$  must orthogonal to past and present observed data .

$$Ee[n]y[n-m] = 0, m \geq 0$$

We want to find  $A_n$  and the  $k_n$  using the above condition.

$e[n]$  is orthogonal to current and past data. First consider the condition that  $e[n]$  is orthogonal to the current data.

$$\begin{aligned} \therefore Ee[n]y[n] &= 0 \\ \Rightarrow Ee[n](x[n] + v[n]) &= 0 \\ \Rightarrow Ee[n]x[n] + Ee[n]v[n] &= 0 \\ \Rightarrow Ee[n](\hat{x}[n] + e[n]) + Ee[n]v[n] &= 0 \\ \Rightarrow Ee^2[n] + Ee[n]v[n] &= 0 \\ \Rightarrow \varepsilon^2[n] + E(x[n] - A_n \hat{x}[n-1] - k_n y[n])v[n] &= 0 \\ \Rightarrow \varepsilon^2[n] - k_n \sigma_v^2 &= 0 \\ \Rightarrow k_n &= \frac{\varepsilon^2[n]}{\sigma_v^2} \end{aligned}$$

We have to estimate  $\varepsilon^2[n]$  at every value of n.

#### **How to do it?**

Consider

$$\begin{aligned} \varepsilon^2[n] &= Ex[n]e[n] \\ &= Ex[n](x[n] - (1 - k_n)a\hat{x}[n-1] - k_n y[n]) \\ &= \sigma_x^2 - (1 - k_n)aEx[n]\hat{x}[n-1] - k_n Ex[n]y[n] \\ &= (1 - k_n)\sigma_x^2 - (1 - k_n)aE(ax[n-1] + w[n])\hat{x}[n-1] \\ &= (1 - k_n)\sigma_x^2 - (1 - k_n)a^2 Ex[n-1]\hat{x}[n-1] \end{aligned}$$

Again

$$\begin{aligned} \varepsilon^2[n-1] &= Ex[n-1]e[n-1] \\ &= Ex[n-1](x[n-1] - \hat{x}[n-1]) \\ &= \sigma_x^2 - Ex[n-1]\hat{x}[n-1] \end{aligned}$$

Therefore,

$$Ex[n-1]\hat{x}[n-1] = \sigma_x^2 - \varepsilon^2[n-1]$$

Hence 
$$\boxed{\varepsilon^2[n] = \frac{\sigma_w^2 + a^2 \varepsilon^2[n-1]}{\sigma_w^2 + \sigma_v^2 + a^2 \varepsilon^2[n-1]} \sigma_v^2}$$

where we have substituted  $\sigma_w^2 = (1-a^2)\sigma_x^2$

We have still to find  $\varepsilon[0]$ . For this assume  $x[-1] = \hat{x}[-1] = 0$ . Hence from the relation

$$\varepsilon[n] = (1-k_n)\sigma_x^2 - (1-k_n)a^2 \varepsilon[n-1]\hat{x}[n-1]$$

we get

$$\varepsilon^2[0] = (1-k_0)\sigma_x^2$$

Substituting we get

$$k_0 = \frac{\varepsilon^2[0]}{\sigma_v^2}$$

in the expression for  $\varepsilon^2[n]$

$$\varepsilon^2[0] = \frac{\sigma_x^2 \sigma_v^2}{\sigma_x^2 + \sigma_v^2}$$

We get

## 8.4 The Scalar Kalman filter algorithm

Given: Signal model parameters  $a$  and  $\sigma_w^2$  and the observation noise variance  $\sigma_v^2$ .

Initialisation  $\hat{x}[-1] = 0$

Step 1.  $n = 0$ . Calculate 
$$\varepsilon^2[0] = \frac{\sigma_x^2 \sigma_v^2}{\sigma_x^2 + \sigma_v^2}$$

Step 2. Calculate 
$$k_n = \frac{\varepsilon^2[n]}{\sigma_v^2}$$

Step 3. Input  $y[n]$ . Estimate  $\hat{x}[n]$  by

$$\hat{x}[n] = a\hat{x}[n-1] + k_n(y[n] - a\hat{x}[n-1])$$

Step 4.  $n = n + 1$ .

Step 5. 
$$\varepsilon^2[n] = \frac{\sigma_w^2 + a^2 \varepsilon^2[n-1]}{\sigma_w^2 + \sigma_v^2 + a^2 \varepsilon^2[n-1]} \sigma_v^2$$

Step 6. Go to Step 2

**Example 2:**

Given

$$x[n] = 0.6x[n-1] + w[n] \quad n \geq 0$$

$$y[n] = x[n] + v[n] \quad n \geq 0$$

$$\sigma_w^2 = 0.25, \sigma_v^2 = 0.5$$

Find the expression for the Kalman filter equations at convergence and the corresponding mean square error.

Using 
$$\varepsilon^2[n] = \frac{\sigma_w^2 + a^2 \varepsilon^2[n-1]}{\sigma_w^2 + \sigma_v^2 + a^2 \varepsilon^2[n-1]} \sigma_v^2$$

We get 
$$\varepsilon^2 = \frac{0.25 + 0.6^2 \varepsilon^2}{0.25 + 0.5 + 0.6 \varepsilon^2} 0.5$$

Solving and taking the positive root

$$\varepsilon^2 = 0.320$$

$$k_n = \varepsilon^2 = 0.390$$

We have to initialize  $\varepsilon^2[0]$ .

Irrespective of this initialization,  $k[n]$  and  $\varepsilon[n]$  converge to final values.

**8.5 Vector Kalman Filter**

$$\mathbf{x}[n] = \mathbf{A}\mathbf{x}[n-1] + \mathbf{w}[n]$$

$$\mathbf{y}[n] = \mathbf{c}\mathbf{x}[n] + \mathbf{v}[n]$$

The vector Kalman filter can be derived in a similar fashion. We will not discuss this derivation.

## **SECTION – IV**

### **SPECTRAL ESTIMATION**

## CHAPTER – 9 : SPECTRAL ESTIMATION TECHNIQUES FOR STATIONARY SIGNALS

### 9.1 Introduction

The aim of spectral analysis is to determine the spectral content of a random process from a finite set of observed data.

- Spectral analysis is a very old problem: Started with the Fourier Series (1807) to solve the wave equation.
- Strum generalized it to arbitrary function (1837)
- Schuster devised periodogram (1897) to determine frequency content numerically.

Consider the definition of the power spectrum of a random sequence  $\{x[n], -\infty < n < \infty\}$

$$S_{XX}(w) = \sum_{m=-\infty}^{\infty} R_{XX}[m] e^{-j2\pi w m}$$

where  $R_{XX}[m]$  is the autocorrelation function.

- The power spectral density is the discrete Fourier Transform (DFT) of the autocorrelation sequence
- The definition involves infinite autocorrelation sequence.
- But we have to use only finite data. This is not only for our inability to handle infinite data, but also for the fact that the assumption of stationarity is valid only for a short duration. For example, the speech signal is stationary for 20 to 80-ms.

Spectral analysis is a preliminary tool – it says that the particular frequency content may be present in the observed signal. Final decision is to be made by ‘Hypothesis Testing’.

Spectral analysis may be broadly classified as *parametric* and *non-parametric*. In the parametric method, the random sequence is modeled by a time-series model, the *model parameters* are estimated from the given data and the spectrum is found by substituting parameters in the model spectrum

We will first discuss the *non-parametric techniques* for spectral estimation. These techniques are based on the Fourier transform of the sample autocorrelation function which is an estimator for the true autocorrelation function.

## 9.2 Sample Autocorrelation Functions

Two estimators of the autocorrelation function exist

$$\hat{R}_{XX}[m] = \frac{1}{N} \sum_{n=0}^{N-1-|m|} x[n]x[n+m] \quad (\text{biased})$$

$$\hat{R}'_{XX}[m] = \frac{1}{N-|m|} \sum_{n=0}^{N-1-|m|} x[n]x[n+m] \quad (\text{unbiased})$$

Note that

$$\begin{aligned} E \hat{R}_{XX}[m] &= \frac{1}{N} \sum_{n=0}^{N-1-|m|} E x[n]x[n+m] \\ &= \frac{N-|m|}{N} R_{XX}[m] \\ &= R_{XX}[m] - \frac{|m|}{N} R_{XX}[m] \\ &= \left( \frac{N-|m|}{N} \right) R_{XX}[m] \end{aligned}$$

Hence  $\hat{R}_{XX}[m]$  is a biased estimator of  $R_{XX}[m]$ . Had we divided the terms under summation by  $N-|m|$  instead of  $N$ , the corresponding estimator would have been unbiased. Therefore,  $\hat{R}'_{XX}[m]$  is an unbiased estimator.

$\hat{R}_{XX}[m]$  is an asymptotically unbiased estimator. As  $N \rightarrow \infty$ , the bias of  $\hat{R}_{XX}[m]$  will tend to 0. The variance of  $\hat{R}_{XX}[m]$  is very difficult to be determined, because it involves fourth-order moments. An approximate expression for the covariance of  $\hat{R}_{XX}[m]$  is given by Jenkins and Watt (1968) as

$$\begin{aligned} \text{Cov}(\hat{R}_{XX}[m_1], \hat{R}_{XX}[m_2]) \\ \cong \frac{1}{N} \sum_{n=-\infty}^{\infty} (R_{XX}[n]R_{XX}[n+m_2-m_1] + R_{XX}[n-m_1]R_{XX}[n+m_2]) \end{aligned}$$

This means that the estimated autocorrelation values are highly correlated.

The variance of  $\hat{R}_{XX}[m]$  is obtained from above as

$$\text{var}(\hat{R}_{XX}[m]) \cong \frac{1}{N} \sum_{n=-\infty}^{\infty} (R_{XX}^2[n] + R_{XX}[n-m]R_{XX}[n+m])$$

Note that the variance of  $\hat{R}_{XX}[m]$  is large for large lag  $m$ , especially as  $m$  approaches  $N$ .



Also as  $N \rightarrow \infty$ ,  $\text{var}(\hat{R}_{XX}[m]) \rightarrow 0$  provided  $\sum_{n=-\infty}^{\infty} (R_{XX}^2[n]) < \infty$ .

As  $N \rightarrow \infty$ ,  $\text{var}(\hat{R}_{XX}[m]) \rightarrow 0$ . Though sample autocorrelation function is a consistent estimator, its Fourier transform is not and here lies the problem of spectral estimation.

Though unbiased and consistent estimators for  $R_{XX}[m]$ ,  $\hat{R}'_{XX}[m]$  is not used for spectral estimation because  $\hat{R}'_{XX}[m]$  does not satisfy the non-negative definiteness criterion.

### 9.3 Periodogram (Schuster, 1898)

$$\hat{S}_{XX}^p(\omega) = \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-j\omega n} \right|^2, \quad -\pi \leq \omega \leq \pi$$

$\hat{S}_{XX}^p(\omega)$  gives the power output of band pass filters of impulse response

$$h_i[n] = \frac{1}{\sqrt{N}} e^{-j\omega_i n} \text{rect}\left(\frac{n}{N}\right)$$

$h_i[n]$  is a very poor band-pass filter.

Also

$$\hat{S}_{XX}^p(\omega) = \sum_{m=-(N-1)}^{N-1} \hat{R}_{XX}[m] e^{-j\omega m}$$

where  $\hat{R}_{XX}[m] = \frac{1}{N} \sum_{n=0}^{N-1-|m|} x[n] x[n+m]$

The periodogram is the Fourier transform of the sample autocorrelation function.

To establish the above relation

Consider

$$\begin{aligned} x_N[n] &= x[n] \text{ for } n < N \\ &= 0 \text{ otherwise} \\ \therefore \hat{R}_{XX}[m] &= \frac{x_N[m] * x_N[-m]}{N} \\ \hat{S}_{XX}^p(\omega) &= \frac{1}{N} \left( \sum_{n=0}^{N-1} x[n] e^{-j\omega n} \right) \left( \sum_{n=0}^{N-1} x[n] e^{j\omega n} \right) \\ &= \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-j\omega n} \right|^2 \end{aligned}$$

$$\begin{aligned}
\text{So } \hat{S}_{XX}^p(w) &= \sum_{m=-(N-1)}^{N-1} \hat{R}_{XX}[m] e^{-jwm} \\
E \hat{S}_{XX}^p(w) &= \sum_{m=-(N-1)}^{N-1} E[\hat{R}_{XX}[m]] e^{-jwm} \\
&= \sum_{m=-(N-1)}^{N-1} \left(1 - \frac{|m|}{N}\right) R_{XX}[m] e^{-jwm}
\end{aligned}$$

as  $N \rightarrow \infty$  the right hand side approaches true power spectral density  $S_{XX}(f)$ . Thus the periodogram is an asymptotically unbiased estimator for the power spectral density.

To prove consistency of the periodogram is a difficult problem.

We consider the simple case when a sequence of *Gaussian white noise* in the following example.

Let us examine the periodogram only at the DFT frequencies  $w_k = \frac{2\pi k}{N}$ ,  $k = 0, 1, \dots, N-1$ .

### **Example 1:**

The periodogram of a zero-mean white Gaussian sequence  $x[n]$ ,  $n = 0, \dots, N-1$ .

The power spectral density is given by

$$S_{XX}(w) = \frac{\sigma_x^2}{2\pi} \quad -\pi < w \leq \pi$$

The periodogram is given by

$$\hat{S}_{XX}^p(w) = \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-jwn} \right|^2$$

Let us examine the periodogram only at the DFT frequencies  $w_k = \frac{2\pi k}{N}$ ,  $k = 0, 1, \dots, N-1$ .

$$\begin{aligned}
\hat{S}_{XX}^p(k) &= \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}kn} \right|^2 \\
&= \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-jw_k n} \right|^2 \\
&= \left( \sum_{n=0}^{N-1} \frac{1}{\sqrt{N}} x[n] \cos w_k n \right)^2 + \left( \sum_{n=0}^{N-1} \frac{x[n] \sin w_k n}{\sqrt{N}} \right)^2 \\
&= C_X^2(w_k) + S_X^2(w_k)
\end{aligned}$$

where  $C_X^2(w_k)$  and  $S_X^2(w_k)$  are the cosine and sine parts of  $\hat{S}_{XX}^p(w)$ .

Let us consider

$$C_X^2(w_K) = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x[n] \cos w_K n$$

which is a linear combination of a Gaussian process.

$$\text{Clearly } E C_X^2(w_K) = 0$$

$$\begin{aligned} \text{var } C_X^2(w_K) &= \frac{1}{N} E \left( \sum_{n=0}^{N-1} x[n] \cos w_K n \right)^2 \\ &= \frac{1}{N} \sum_{n=0}^{N-1} E x^2[n] \cos^2 w_K n + E(\text{Cross terms}) \end{aligned}$$

Assume the sequence  $x[n]$  to be independent.

Therefore,

$$\begin{aligned} \text{var}(C_K(w_K)) &= \frac{\sigma_X^2}{N} \sum_{n=0}^{N-1} \cos^2 w_K n = \frac{\sigma_X^2}{N} \sum_{n=0}^{N-1} \left( \frac{1 + \cos 2w_K n}{2} \right) \\ &= \frac{\sigma_X^2}{N} \left( \frac{N}{2} + \cos(N-1)w_K \frac{\sin Nw_K}{2 \sin w_K} \right) \\ &= \sigma_X^2 \left( \frac{1}{2} + \cos(N-1)w_K \frac{\sin Nw_K}{2N \sin w_K} \right) \end{aligned}$$

$$\text{For } w_K = \frac{2\pi}{N} k$$

$$\frac{\sin Nw_K}{2 \sin w_K} = 0 \quad k \neq 0, k \neq \frac{N}{2} \text{ (assuming } N \text{ even).}$$

$$\therefore \text{var}(C_K(w_K)) = \frac{\sigma_X^2}{2} \text{ for } k \neq 0$$

$$\text{Again } \frac{\sin Nw_K}{N \sin w_K} = 1 \text{ for } k = 0.$$

$$\therefore \text{var}(C_K(w_K)) = \sigma_X^2 \text{ for } k = 0, k = \frac{N}{2}$$

Similarly considering the sine part

$$S_K(w_K) = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x[n] \sin w_K n$$

$$E S_K(w_K) = 0 \quad \because x[n] \text{ is zero mean.}$$

$$\text{var } S_K(w_K) = 0 \text{ for } k = 0 \text{ (for dc part has no sine term)}$$

$$= \frac{\sigma_X^2}{2} \text{ for } k \neq 0.$$

Therefore for  $k = 0$ , Distribution.

$$C_K[k] \sim N(0, \sigma_X^2)$$

$$S_K[k] = 0.$$

for  $k \neq 0$

$$C_k \sim N\left(0, \frac{\sigma_x^2}{2}\right)$$

$$S_k \sim N\left(0, \frac{\sigma_x^2}{2}\right).$$

We can also show that

$$\text{Cov}(S_X(w_k)C_X(w_k)) = 0.$$

So  $C_X(w_k)$  and  $S_X(w_k)$  are independent Gaussian R.V.s..

The periodogram

$$\hat{S}_{XX}^p[k] = C_X^2[k] + S_X^2[k] \quad \text{has a chi-square distribution}$$

## 9.4 Chi square distribution

Let  $X_1, X_2, \dots, X_N$  be independent zero-mean Gaussian variables each with variance  $\sigma_X^2$ . Then  $Y = X_1^2 + X_2^2 + \dots + X_N^2$  has (chi square)  $\chi_N^2$  distribution with mean

$$EY = E(X_1^2 + X_2^2 + \dots + X_N^2) = N\sigma_X^2 \quad \text{and variance } 2N\sigma_X^4.$$

$$\hat{S}_{XX}[k] = C_X[k] + S_X^2[k].$$

It is a  $\chi_2^2$  distribution.

$$E \hat{S}_{XX}[k] = \frac{\sigma_X^2}{2} + \frac{\sigma_X^2}{2} = \sigma_X^2 = S_{XX}[k]$$

$\Rightarrow \hat{S}_{XX}[k]$  is unbiased

$$\text{var}(\hat{S}_{XX}[k]) = 2 \times 2 \left( \frac{\sigma_X^2}{2} \right)^2 = S_{XX}^2[k] \quad \text{which is independent of } N$$

$\hat{S}_{XX}[0] = C_0^2[0]$  is a  $\chi_1^2$  distribution of degree of freedom 1

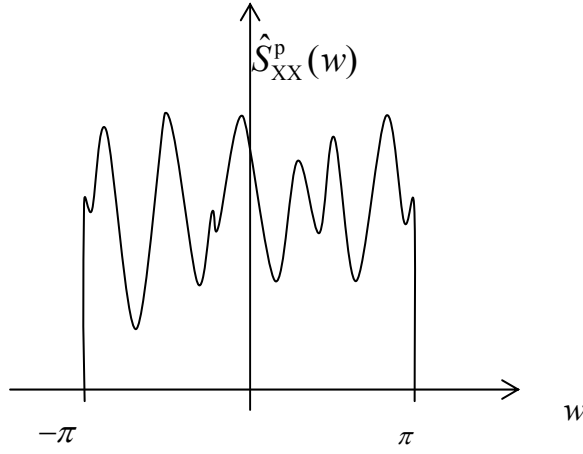
$$E \hat{S}_{XX}[0] = \sigma_x^2$$

$\Rightarrow \hat{S}_{XX}[0]$  is unbiased

$$\text{and var}(\hat{S}_{XX}[0]) = \sigma_X^4 = S_{XX}^2[0].$$

It can be shown that for the Gaussian independent white noise sequence at any frequency  $w$ ,

$$\text{var}(\hat{S}_{XX}^p(w)) = S_{XX}^2(w)$$



### For general case

$$\hat{S}_{XX}^p(\omega) = \sum_{m=-(N-1)}^{N-1} \hat{R}_{XX}[m] e^{-j\omega m}$$

where  $\hat{R}_{XX}[m] = \frac{1}{N} \sum_{n=0}^{N-1-|m|} x[n]x[n+m]$ , the biased estimator of autocorrelation fn.

$$= \sum_{m=-(N-1)}^{N-1} \left(1 - \frac{|m|}{N}\right) \hat{R}'_{XX}[m] e^{-j\omega m}$$

$$= \sum_{m=-(N-1)}^{N-1} w_B[m] \hat{R}'_{XX}[m] e^{-j\omega m}$$

where  $\hat{R}'_{XX}[m]$  is the unbiased estimator of auto correlation

and  $w_B[m] = 1 - \frac{|m|}{N}$  is the Bartlett Window.

( Fourier Transform of product of two functions)

So,  $\hat{S}_{XX}^p(\omega) = W_B(\omega) * FT\{\hat{R}'_{XX}[m]\}$

$$\begin{aligned} E \hat{S}_{XX}^p(\omega) &= W_B(\omega) * FT\{E \hat{R}'_{XX}[m]\} \\ &= W_B(\omega) * S_{XX}(\omega) \\ &= \int_{-\pi}^{\pi} W_B(\omega - \xi) S_{XX}(\xi) d\xi \end{aligned}$$

As  $N \rightarrow \infty$ ,  $E \hat{S}_{XX}^p(\omega) \rightarrow S_{XX}(\omega)$

Now  $\text{var}(\hat{S}_{XX}^p(\omega))$  cannot be found out exactly (no analytical tool). But an approximate expression for the covariance is given by

$$\text{Cov}(\hat{S}_{XX}^p(\omega_1), \hat{S}_{XX}^p(\omega_2))$$

$$\simeq S_{XX}(\omega_1) S_{XX}(\omega_2) \left[ \left( \frac{\sin \frac{N(\omega_1 + \omega_2)}{2}}{N \sin(\frac{\omega_1 + \omega_2}{2})} \right)^2 + \left( \frac{\sin \frac{N(\omega_1 - \omega_2)}{2}}{N \sin(\frac{\omega_1 - \omega_2}{2})} \right)^2 \right]$$

$$\text{var}(\hat{S}_{XX}^p(w)) \simeq S_{XX}^2(w) \left[ 1 + \left( \frac{\sin Nw}{N \sin w} \right)^2 \right]$$

$$\therefore \text{var} \{ \hat{S}_{XX}^p(w) \} \simeq 2S_{XX}^2(w) \quad \text{for } w = 0, \pi$$

$$\cong S_{XX}^2(w) \quad \text{for } 0 < w < \pi$$

Consider  $w_1 = 2\pi \frac{k_1}{N}$  and  $w_2 = 2\pi \frac{k_2}{N}$ ,  $k_1, k_2$  integers.

Then

$$\text{Cov}(\hat{S}_{XX}^p(w_1), \hat{S}_{XX}^p(w_2)) \cong 0$$

This means that there will be no correlation between two neighboring spectral estimates.

Therefore periodogram is not a reliable estimator for the power spectrum for the following two reasons:

- (1) The periodogram is not a consistent estimator in the sense that  $\text{var}(\hat{S}_{XX}^p(w))$  does not tend to zero as the data length approaches infinity.
- (2) For two frequencies, the covariance of the periodograms decreases as data length increases. Thus the periodogram is erratic and widely fluctuating.

Our aim will be to look for spectral estimator with variance as low as possible, and without increasing much the bias.

## 9.5 Modified Periodograms

### Data windowing

Multiply *data* with a *more suitable window* rather *than* rectangular before finding the periodogram. The resulting periodogram will have less bias.

### 9.5.1 Averaged Periodogram: The Bartlett Method

The motivation for this method comes from the observation that

$$\lim_{N \rightarrow \infty} E \hat{S}_{XX}^p(w) \rightarrow S_{XX}(w)$$

We have to modify  $\hat{S}_{XX}^p(w)$  to get a consistent estimator for  $S_{XX}(f)$ .

Given the data  $x[n], n = 0, 1, \dots, N-1$

Divide the data in  $K$  non-overlapping segments each of length  $L$ .

Determine periodogram of each section.

$$\hat{S}_{XX}^{(k)}(w) = \frac{1}{L} \left| \sum_{n=0}^{L-1} x[n] e^{-jwn} \right|^2 \quad \text{for each section}$$

$$\text{Then } \hat{S}_{XX}^{(av)}(w) = \frac{1}{K} \sum_{m=0}^{K-1} \hat{S}_{XX}^{(k)}(w).$$

As shown earlier,

$$E\hat{S}_{XX}^{(k)}(w) = \int_{-\pi}^{\pi} W_B(w - \xi) S_{XX}(\xi) d\xi$$

$$\text{where } w_B[m] = \begin{cases} 1 - \frac{|m|}{L}, & |m| \leq M-1 \\ 0, & \text{otherwise} \end{cases}$$

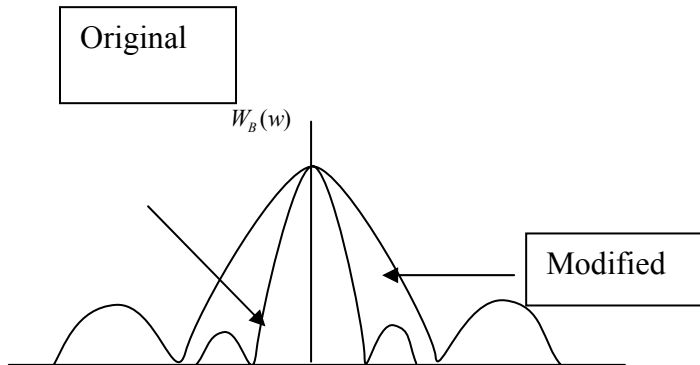
$$W_B(w) = \left( \frac{\sin \frac{wL}{2}}{\sin \frac{w}{2}} \right)^2$$

$$\hat{S}_{XX}^{(k)}(w) = \sum_{m=-(L-1)}^{L-1} \hat{R}_{XX}[m] e^{-jwm}$$

$$E\hat{S}_{XX}^{av}(w) = \frac{1}{K} \sum_{k=0}^{K-1} E\hat{S}_{XX}^{(k)}(w) = E\{S_{XX}^{(k)}(w)\}$$

$$E\hat{S}_{XX}^{av}(w) = \int_{-\pi}^{\pi} W_B(w - \xi) S_{XX}(\xi) d\xi.$$

To find the mean of the averaged periodogram, the true spectrum is now convolved with the frequency  $W_B(f)$  of the Barlett window. The effect of reducing the length of the data from  $N$  points to  $L = N/K$  results in a window whose spectral width has been increased by a factor  $K$  consequently the frequency resolution has been reduced by a factor  $K$ .



**Figure Effect of reducing window size**

### 9.5.2 Variance of the averaged periodogram

$\text{var}(\hat{S}_{XX}^{av}(w))$  is not simple to evaluate as 4<sup>th</sup> order moments are involved.

#### Simplification :

Assume the  $K$  data segments are independent.

$$\begin{aligned}\text{Then } \text{var}(\hat{S}_{XX}^{av}(w)) &= \text{var}\left\{\sum_{k=0}^{K-1} \hat{S}_{XX}^{(k)}(w)\right\} \\ &= \frac{1}{K^2} \sum_{k=0}^{K-1} \text{var} \hat{S}_{XX}^{(m)}(w) \\ &\simeq \frac{1}{K^2} \times K \left(1 + \left(\frac{\sin wL}{L \sin w}\right)^2\right) S_{XX}^2(w) \\ &\simeq \frac{1}{K} \text{original variance of the periodogram.}\end{aligned}$$

So variance will be reduced by a factor less than  $K$  because in practical situations, data segments will not be independent.

For large  $L$  and large  $K$ ,  $\text{var}(\hat{S}_{XX}^{av}(w))$  will tend to zero and  $\hat{S}_{XX}^{av}(w)$  will be a consistent estimator.

#### The Welch Method (Averaging modified periodograms)

- (1) Divide the data into overlapping segments ( overlapping by about 50 to 75%).
- (2) Window the data so that the *modified periodogram* of each segment is

$$\begin{aligned}\hat{S}_{XX}^{(\text{mod})}(w) &= \frac{1}{UL} \left| \sum_{n=0}^{L-1} x[n]w[n]e^{-jwn} \right|^2 \\ &= \frac{1}{L} \sum_{n=0}^{L-1} w^2[n]\end{aligned}$$

The window  $w[n]$  need not be an even function and is used to control spectral leakage.

$$\begin{aligned}\sum_{n=0}^{L-1} x[n]w[n]e^{-jwn} &\text{ is the DTFT of } x[n]w[n] \text{ where } w[n] = 1 \text{ for } n = 0, \dots, L-1 \\ &= 0 \text{ otherwise}\end{aligned}$$

- (3) Compute  $\hat{S}_{XX}^{(Welch)}(w) = \frac{1}{K} \sum_{k=0}^{K-1} \hat{S}_{XX}^{(\text{mod})}(w)$ .



## 9.6 Smoothing the periodogram : The Blackman and Tukey Method

- Widely used non parametric method
- Biased autocorrelation sequence is used  $\hat{R}_{XX}[m]$ .
- Higher order or large lags autocorrelation estimation involves less data ( $N - m$ ,  $m$  large) and so more prone to error, we give less importance to higher-order autocorrelation.

$\hat{R}_{XX}[m]$  is multiplied by a window sequence  $w[m]$  which under weighs the autocorrelation function at large lags. The window function  $w[m]$  has the following properties.

$$\begin{aligned} 0 < w[m] &< 1 \\ w[0] &= 1 \\ w[-m] &= w[m] \\ w[m] &= 0 \text{ for } |m| > M. \end{aligned}$$

$w[0] = 1$  is a consequence of the fact that the smoothed periodogram should not modify a smooth spectrum and so  $\int_{-\pi}^{\pi} W(w)dw = 1$ .

The smoothed periodogram is given by

$$\hat{S}_{XX}^{BT}(w) = \sum_{m=-(M-1)}^{M-1} w[m] \hat{R}_{XX}[m] e^{-jwm}$$

### Issues concerned :

#### 1. How to select $w[m]$ ?

There are a large numbers of windows available. Use a window with small side-lobes. This will reduce bias and improve resolution.

#### 2. How to select $M$ . – Normally

$$M \simeq \frac{N}{5} \text{ or } M \simeq 2\sqrt{N} \text{ (Mostly based on experience)}$$

$$N \sim 1000 \quad \text{if } N \text{ is large } 10,000$$

$\hat{S}_{XX}^{BT}(w)$  = convolution of  $\hat{S}_{XX}^p(w)$  and  $W(w)$ , the F.T. of the window sequence.  
= Smoothing the periodogram, thus decreasing the variance in the estimate at the expense of reducing the resolution.

$$E(\hat{S}_{XX}^{BT}(w)) = E \hat{S}_{XX}^p(w) * W(w)$$

$$\text{where } E \hat{S}_{XX}^p(w) = \int_{-\pi}^{\pi} S_{XX}(\theta) W_B(w - \theta) d\theta$$

or from time domain

$$\begin{aligned} E \hat{S}_{XX}^{BT}(w) &= E \sum_{m=-(M-1)}^{M-1} W[m] \hat{R}_{XX}[m] e^{-jwm} \quad \text{asymptotically unbiased} \\ &= \sum_{m=-(M-1)}^{M-1} W[m] \underbrace{E \hat{R}_{XX}[m]} e^{-jwm} \end{aligned}$$

$\hat{S}_{XX}^{BT}(w)$  can be proved to be asymptotically unbiased.

$$\text{and variance of } \hat{S}_{XX}^{BT}(w) \simeq \frac{S_{XX}^2(w)}{N} \sum_{K=-(M-1)}^{M-1} w^2[k]$$

Some of the popular windows are rectangular window, Bartlett window, Hamming window, Hanning window, etc.

### **Procedure**

1. Given the data  $x[n], n = 0, 1, \dots, N-1$
2. Find periodogram.  $\hat{S}_{XX}^p(2\pi k / N) = \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-j2\pi k n / N} \right|^2$ ,
3. By IDFT find the autocorrelation sequence.
4. Multiply by proper window and take FFT.

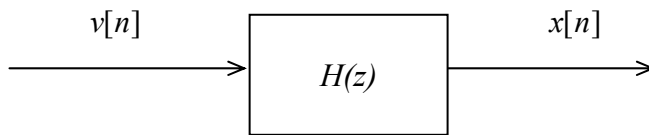
## **9.7 Parametric Method**

Disadvantage classical spectral estimators like Blackman Tuckey Method using windowed autocorrelation function.

- Do not use a priori information about the spectral shape
- Do not make realistic assumptions about  $x[n]$  for  $n < 0$  and  $n \geq N$ .

Normally  $\hat{R}_X[m], m = 0, \pm 1, \dots, \pm(N-1)$  can be estimated from  $N$  sample values. From these autocorrelations we can estimate a model for the signal – basically a time series model. Once the model is available it is equivalent to know the autocorrelation for all lags and hence will give better spectral estimation. (better resolution) A stationary signal can be represented by ARMA (p, q) model.

$$x[n] = \sum_{i=1}^p a_i x[n-i] + \sum_{i=0}^q b_i v[n-i]$$



$$H(z) = \frac{B(z)}{A(z)} = \frac{\sum_{i=0}^q b_i z^{-i}}{1 - \sum_{i=1}^p a_i z^{-i}}$$

and

$$S_{xx}(w) = |H(w)|^2 \sigma_v^2$$

- Model signal as AR/MA/ARMA process
- Estimate *model parameters* from the given data
- Find the power spectrum by substituting the values of the model parameters in expression of power spectrum of the model

## 9.8 AR spectral estimation

$AR(p)$  process is the most popular model based technique.

- Widely used for parametric spectral analysis because:
- Can approximate continuous power spectrum arbitrarily well
- (but might need very large  $p$  to do so)
- Efficient algorithms available for model parameter estimation
- if the process Gaussian maximum entropy spectral estimate is  $AR(p)$  spectral estimate
- Many physical models suggest AR processes (e.g. speech)
- Sinusoids can be expressed as AR-like models

The spectrum is given by

$$\hat{S}_{xx}(w) = \frac{\sigma_v^2}{\left| 1 - \sum_{i=1}^p a_i e^{-jwi} \right|^2}$$

where  $a_i$ s are  $AR(p)$  process parameters.

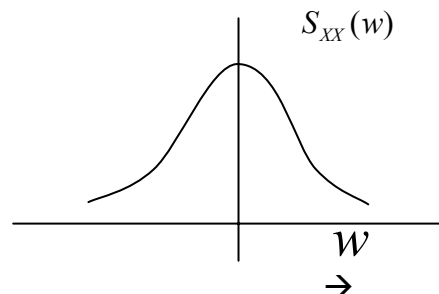


Figure an AR spectrum

## 9.9 The Autocorrelation method

The  $a_i$ s are obtained by solving the Yule Walker equations corresponding to estimated autocorrelation functions.

$$\begin{bmatrix} \hat{R}_{xx}[0] & \hat{R}_{xx}[1] & \dots & \hat{R}_{xx}[p-1] \\ \hat{R}_{xx}[1] & \hat{R}_{xx}[0] & \dots & \hat{R}_{xx}[p-2] \\ \dots & \dots & \dots & \dots \\ \hat{R}_{xx}[p-1] & \hat{R}_{xx}[1] & \dots & \hat{R}_{xx}[0] \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \dots \\ a_p \end{bmatrix} = \begin{bmatrix} \hat{R}_x[1] \\ \hat{R}_x[2] \\ \dots \\ \hat{R}_x[p+1] \end{bmatrix}$$

$$\sigma_v^2 = \hat{R}_x[0] - \sum_{i=1}^p a_i \hat{R}_{xx}[i]$$

## 9.10 The Covariance method

The problem of estimating the AR-parameters may be considered in terms of minimizing the sum-square error

$$\varepsilon[n] = \sum_{n=p}^{N-1} e^2[n]$$

with respect to the AR parameters, where

$$e[n] = x[n] - \sum_{i=1}^p a_i x[n-i]$$

This formulation results in

$$\begin{bmatrix} \hat{R}_{xx}[1,1] & \hat{R}_{xx}[2,1] & \dots & \hat{R}_{xx}[p,1] \\ \hat{R}_{xx}[1,2] & \hat{R}_{xx}[2,2] & \dots & \hat{R}_{xx}[p,2] \\ \dots & \dots & \dots & \dots \\ \hat{R}_{xx}[1,p] & \hat{R}_{xx}[2,p] & \dots & \hat{R}_{xx}[p,p] \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \dots \\ a_p \end{bmatrix} = \begin{bmatrix} \hat{R}_x[0,1] \\ \hat{R}_x[0,2] \\ \dots \\ \hat{R}_x[0,p+1] \end{bmatrix}$$

$$\sigma_v^2 = \hat{R}_x[0,0] - \sum_{i=1}^p a_i \hat{R}_{xx}[0,i]$$

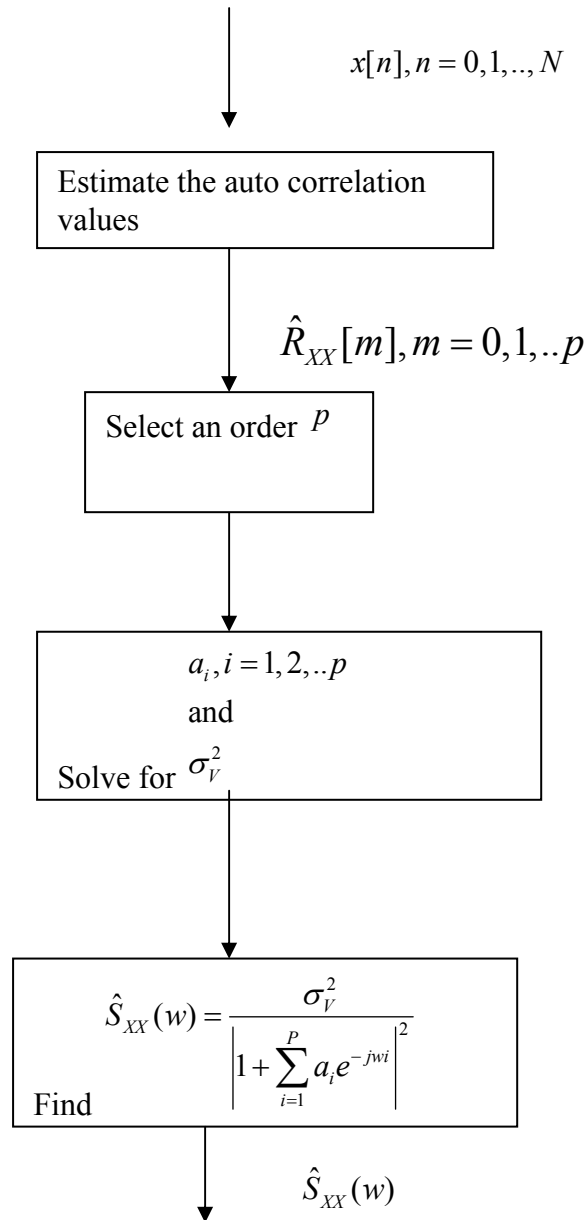
where

$$\hat{R}_{xx}[k,l] = \frac{1}{N-p} \sum_{n=p}^N x[n-k]x[n-l]$$

is an estimate for the autocorrelation function.

Note that the autocorrelation matrix in the *Covariance method* is not Toeplitz and cannot be solved by efficient algorithms like the Levinson Durbin recursion.

The flow chart for the AR spectral estimation is given below:



### Some questions :

*Can an  $AR(p)$  process model a band pass signal?*

- If we use  $AR(1)$  model, it will never be able to model a band-pass process. If one sinusoid is present then  $AR(2)$  process will be able to discern it. If there is a strong frequency component  $w_0$ , an  $AR(2)$  process with poles at  $re^{\pm jw_0}$  with  $r \rightarrow 1$  will be able to discriminate the frequency components.

*How do we select the order of the  $AR(p)$  process?*

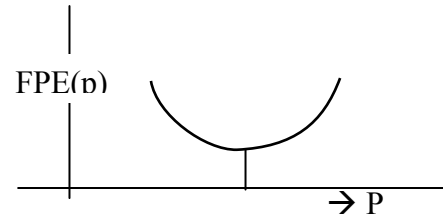
- MSE will give some guidance regarding the selected order is proper or not.  
For spectral estimation, some criterion function with respect to the order parameter  $p$  are to be minimized. For example,

- Minimize Forward Prediction Error.

$$FPE(p) = \hat{\sigma}_p^2 \frac{N + p + 1}{N - p - 1}$$

where N = No of data

$\hat{\sigma}_p^2$  = mean square prediction error (variance for non zero mean case)



- Akaike Information Criteria

-Minimize

$$AIC(p) = \ln(\hat{\sigma}_v^2) + \frac{2p}{N}$$

### 9.11 Frequency Estimation of Harmonic signals

A class of spectral estimators is based on eigen decomposition of the autocorrelation matrix of the data vector. Either the eigen values or the eigen vectors can be used. Notable of these algorithms is the MUSIC (Multiple Signal Classification) algorithm. We will discuss this algorithm also.

## **10. Text and Reference**

1. **M. Hays, *Statistical Digital Signal Processing and Modelling*, John Willey and Sons, 1996.**
2. M.D. Srinath, P.K. Rajasekaran and R. Viswanathan, *Statistical Signal Processing with Applications*, PHI, 1996.
3. Simon Haykin, *Adaptive Filter Theory*, Prentice Hall, 1996
4. D.G. Manolakis, V.K. Ingle and S.M. Kogon, *Statistical and Adaptive Signal Processing*, McGraw Hill, 2000
5. S. M. Kay, *Modern Spectral Estimation*, Prentice Hall, 1987
6. S. J. Orfanidis, *Optimum Signal Processing*, Second Edition, MacMillan Publishing, 1989.
7. H. Stark and J.W. Woods, *Probability and Random Processes with Applications to Signal Processing*, Prentice Hall 2002.
8. A. Papoulis and S.U. Pillai, *Probability, Random Variables and Stochastic Processes*, 4th Edition, McGraw-Hill, 2002