

LECTURE NOTES FOR THE COURSE SC42150

---

# Statistical Signal Processing

Stochastic Processes for Scientists and Engineers with Modern Applications

---

**Carlos S. Smith**  
**Michel Verhaegen**

---

September 23, 2020

Delft University of Technology  
Delft Center for Systems and Control  
*Control for Scientific Imaging (CSI) Group*  
Mekelweg 2, NL-2628 CD Delft, The Netherlands  
M.Verhaegen@tudelft.nl  
C.S.Smith@tudelft.nl



# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>1</b>  |
| 1.1      | Empirical or data-driven modeling . . . . .                       | 1         |
| 1.2      | Increasing relevance of stochastic processes in Physics . . . . . | 2         |
| 1.3      | In the footsteps of Carl Friedrich Gauss . . . . .                | 5         |
| 1.4      | Objectives of the Course . . . . .                                | 9         |
| 1.4.1    | Objective 1: Description of stochastic processes . . . . .        | 9         |
| 1.4.2    | Objective 2: Generation of stochastic processes . . . . .         | 9         |
| 1.4.3    | Objective 3: Parameter Estimation . . . . .                       | 10        |
| 1.4.4    | Objective 4: Optimal Filtering . . . . .                          | 10        |
| 1.5      | Applications of Optimal Filtering . . . . .                       | 12        |
| 1.5.1    | Denoising . . . . .   | 12        |
| 1.5.2    | Deconvolution . . . . .   | 12        |
| 1.5.3    | Prediction . . . . .  | 13        |
| 1.5.4    | Active Noise Cancellation . . . . .                               | 13        |
|          | Exercises . . . . .   | 15        |
| <b>2</b> | <b>Signals and Systems</b>  | <b>17</b> |
| 2.1      | Introduction . . . . .  | 18        |
| 2.2      | Discrete-Time Signals . . . . .                                   | 18        |
| 2.2.1    | Definition and examples . . . . .                                 | 18        |
| 2.2.2    | The Discrete-Time Fourier Transform . . . . .                     | 19        |
| 2.2.3    | The z-Transform . . . . .   | 21        |
| 2.3      | Discrete-Time Systems . . . . .                                   | 22        |
| 2.3.1    | Definition . . . . .  | 22        |
| 2.3.2    | Linear Time-invariant (LTI) Discrete Time Systems . . . . .       | 23        |
| 2.3.3    | Stability, Causality and Minimum-phase LTI systems . . . . .      | 24        |
| 2.4      | Optimizing cost functions w.r.t. a complex parameter . . . . .    | 28        |
|          | Exercises . . . . .   | 30        |
| <b>3</b> | <b>Random Variables</b>   | <b>33</b> |
| 3.1      | Introduction . . . . .  | 34        |
| 3.2      | Random Variables . . . . .  | 34        |
| 3.2.1    | Discrete Random Variables . . . . .                               | 35        |
| 3.2.2    | Continuous Random Variables . . . . .                             | 36        |
| 3.2.3    | Characterizing Random Variables . . . . .                         | 37        |
| 3.2.4    | Two Random Variables . . . . .                                    | 40        |

|          |  |            |
|----------|--|------------|
| 3.2.5    | Gaussian Random Variables . . . . .                                | 42         |
|          | Exercises . . . . .  | 43         |
| <b>4</b> | <b>Estimation</b>  | <b>47</b>  |
| 4.1      | Introduction . . . . .   | 48         |
| 4.1.1    | Basic definitions . . . . .  | 48         |
| 4.2      | Linear regression . . . . .  | 52         |
| 4.2.1    | Introduction . . . . .   | 52         |
| 4.2.2    | Linear least-squares estimation . . . . .                          | 53         |
| 4.2.3    | Linear regression as a statistical estimation problem . . . . .    | 56         |
| 4.2.4    | Weighted linear least squares estimation . . . . .                 | 61         |
| 4.3      | The Cramér-Rao lower bound . . . . .                               | 64         |
| 4.4      | Maximum likelihood estimator . . . . .                             | 69         |
|          | Exercises . . . . .  | 77         |
| <b>5</b> | <b>Stochastic Processes</b>  | <b>83</b>  |
| 5.1      | Introduction . . . . .   | 84         |
| 5.2      | Stochastic Processes . . . . .                                     | 84         |
| 5.2.1    | Ensemble Averages for a Stochastic Process . . . . .               | 86         |
| 5.2.2    | Ensemble Averages for two Stochastic Processes . . . . .           | 87         |
| 5.2.3    | Gaussian Processes . . . . .                                       | 87         |
| 5.2.4    | Stationary Processes . . . . .                                     | 88         |
| 5.2.5    | Wide Sense Stationary (WSS) Processes . . . . .                    | 89         |
| 5.2.6    | Autocorrelation matrix of a Stochastic Process . . . . .           | 91         |
| 5.2.7    | Ergodicity . . . . .   | 91         |
| 5.2.8    | WSS Processes in the Frequency Domain . . . . .                    | 94         |
|          | Exercises . . . . .  | 99         |
| <b>6</b> | <b>Filtering Stochastic Processes</b>                              | <b>103</b> |
| 6.1      | Introduction . . . . .   | 104        |
| 6.2      | General mixed causal, anti-causal LTI systems . . . . .            | 104        |
| 6.2.1    | The Auto- and Cross correlation Function after filtering . . . . . | 106        |
| 6.2.2    | The Power and Cross Spectrum after filtering . . . . .             | 107        |
| 6.3      | Zero-Mean White Noise . . . . .                                    | 108        |
| 6.4      | ARMA, AR, MA models . . . . .                                      | 110        |
| 6.4.1    | Definition of the models . . . . .                                 | 110        |
| 6.4.2    | Calculation of the Power Spectrum . . . . .                        | 112        |
| 6.4.3    | Calculation of the Auto-Correlation Function . . . . .             | 115        |
|          | Exercises . . . . .  | 119        |
| <b>7</b> | <b>Inverse Problems in Time and Frequency Domain</b>               | <b>125</b> |
| 7.1      | Introduction . . . . .   | 126        |
| 7.2      | A Motivating Example . . . . .                                     | 127        |
| 7.3      | Spectral Factorization . . . . .                                   | 128        |
| 7.3.1    | Problem Definition . . . . .                                       | 128        |
| 7.3.2    | Solution . . . . .   | 129        |
| 7.3.3    | Use of the Spectral Factorization . . . . .                        | 133        |

|          |  |            |
|----------|--|------------|
| 7.4      | Finding the Shaping filter of a stochastic process given its Auto-correlation Function . . . . . | 134        |
| 7.4.1    | Problem Formulation . . . . .  | 135        |
| 7.4.2    | Solution . . . . .   | 136        |
|          | Exercises . . . . .  | 143        |
| <b>8</b> | <b>Parameter Estimation: The Linear Least Squares Method</b>                                     | <b>147</b> |
| 8.1      | Introduction . . . . .   | 148        |
| 8.2      | The Linear Least Squares (LLSQ) Problem . . . . .  | 149        |
| 8.2.1    | The LLSQ problem for estimating the parameters of an AR model . . . . .                          | 149        |
| 8.2.2    | Linear Regression . . . . .  | 152        |
| 8.3      | Solution to the AR parameter estimation problem . . . . .  | 154        |
| 8.3.1    | Deriving the solution . . . . .  | 154        |
| 8.3.2    | The Orthogonality Condition . . . . .  | 157        |
| 8.4      | The accuracy of the LLSQ estimated parameters . . . . .  | 158        |
|          | Exercises . . . . .  | 160        |
| <b>9</b> | <b>Optimal Filtering</b>   | <b>165</b> |
| 9.1      | Introduction . . . . .   | 166        |
| 9.2      | Generic Filtering Problem . . . . .  | 166        |
| 9.3      | Minimum Variance FIR Wiener Filter . . . . .   | 169        |
| 9.3.1    | The Generic Problem Formulation . . . . .  | 169        |
| 9.3.2    | Solution to the Generic Problem . . . . .  | 169        |
| 9.3.3    | Application to the Denoising Problem . . . . .   | 172        |
| 9.3.4    | Application to the Prediction Problem . . . . .  | 173        |
| 9.3.5    | Application to the Active Noise Cancellation Problem . . .                                       | 174        |
| 9.4      | Minimum Variance IIR Wiener Filter . . . . .   | 175        |
| 9.4.1    | The Generic Problem Formulation . . . . .  | 175        |
| 9.4.2    | The mixed causal, anti-causal IIR Wiener Filter . . . . .  | 176        |
| 9.4.3    | Application of the mixed causal, anti-causal IIR Wiener filter to Denoising . . . . .            | 177        |
| 9.4.4    | Application of the mixed causal, anti-causal IIR Wiener filter to Deconvolution . . . . .        | 178        |
| 9.4.5    | The Causal IIR Wiener Filter . . . . .   | 179        |
| 9.5      | Example . . . . .  | 181        |
| 9.5.1    | Calculating the mixed causal, anti-causal IIR Wiener filter                                      | 182        |
| 9.5.2    | Calculating the causal IIR Wiener filter . . . . .   | 183        |
|          | Exercises . . . . .  | 184        |
|          | <b>Subject Index</b>   | <b>189</b> |



# Chapter 1

## Introduction

---

### 1.1 Empirical or data-driven modeling

A key driver in physical research is to understand and harness physical reality around us and inside us. This understanding usually comes in the form of explaining physical phenomena by laws of nature. Examples of such laws are Huygens-Fresnel principles to describe the propagation of optical wave fronts. These laws come with assumptions and idealistic boundary conditions. Thus they are mathematical models that described an idealized world, and only in part.

A crucial role in deducing these laws is played by the careful and critical observation of reality. This was in ancient times the observation through our senses but in modern times these observations constitute of carefully designed and often complicated and expensive experiments. Experimentation is used in various ways. It can be used to discover new laws of nature, establish universal constants or to verify postulated hypotheses to describe physical phenomena. In these course notes such physical phenomena will be called physical systems or systems in brief.

The interaction between measurements from dedicated instruments (often designed and operated by the scientist himself) and mathematics to derive or refine physical models has for centuries been a non-separable marriage. As the instruments are the eyes of the scientist, the interpretation and processing of the measurements are often calling upon profound mathematical skills.

Generally the branch of science dedicated to deriving physical laws or mathematical descriptions to relate measurements of physical quantities is called *empirical modeling* or *data driven modeling*. Empirical or data driven modeling covers a vast variety of modeling approaches that employ measurements to derive mathematical models to approximate the behavior of systems, going from calibrating unknown parameters in mathematical model structures derived from first principles to deriving full mathematical relationships directly from measurements. An example of calibrating mathematical models is the use of mathematical models describing diffraction of light to determine the wavelength of a coherent laser source (<http://academia.hixie.ch/bath/laser/home.html>).

An example of deriving full mathematical relationships is the modeling of the dynamic relationship between the voltage (input) and current (output) of an unknown electric circuit in a box. By applying a chosen voltage signal to the circuit, its corresponding output is measured. Both signals are then used to derive a mathematical relationship (such as a differential or difference equation) that has the current as its solution for the given voltage as input.

Key in addressing data driven modeling problems is the ability to cope with uncertainty, both in the measurement devices due to their finite precision as well as in the discrepancy between the behavior predicted by the mathematical model and the real life measurements.

Generally in reality the relationships that we are considering are dynamic, e.g. difference equations. Therefore the uncertainty to be considered is not that of a single number or a limited number of quantities, but of time series.

The goal of this course is to provide an introduction towards the description and use of (time or space) sequences of “errors” using statistics.

We will start in Chapter 3 with a brief review on the description of a finite number of random variables and extend these notions to infinite sequences of random variables. Such sequences will be called *random* or *stochastic* processes. For this introductory course we will mainly focus on stationary random processes that are solely described by their so-called second order statistical moments, that is their auto- or cross-correlation functions in the time or spatial domain or power or cross-spectra in the frequency domain.

In the following sections of this introductory chapter we briefly outline a number of key discoveries that have greatly influenced the field of empirical modeling.

---

## 1.2 Increasing relevance of stochastic processes in Physics

The incentives characterizing modern physics to study and master systems at an increasingly smaller dimension, naturally lead to an increase in the importance of unpredictable signals. Random phenomena that could at larger dimensions be neglected suddenly may become even dominant at a small scale.

A few examples are:

1. The collision of molecules in a liquid or air substance with other masses may create a dynamic motion of these masses, that was first discovered by Robert Brown in 1827 [1]. He used a simple microscope to study the movement of particles from pollen in water immersion. He observed that pollen particles of the size of  $5\mu m$  immersed in liquid were not at rest.
2. Thermal effects induces resistor noise that may become dominant in electric measuring devices when measuring small distances.
3. Turbulence induces a variation in the refractive index and as a consequence variation in the optical path length through a turbulent atmosphere. This



phenomenon causes a 30 meter ground based telescope to have a similar resolution compared to one of the first 6 inch ( $\approx 15\text{cm}$ ) reflecting telescopes designed under the impulse of Isaac Newton in the 17-th century.

4. Quantum mechanics shows that the behavior of quantum particles is described by probabilistics. The effect of quantum tunneling is but one example that explains the probability of an electron tunneling through a thin wall. It is a very sensitive phenomena that forms the basis of many highly sensitive measuring techniques, such as as in scanning tunneling microscopy.

A first more elaborated example is given next.

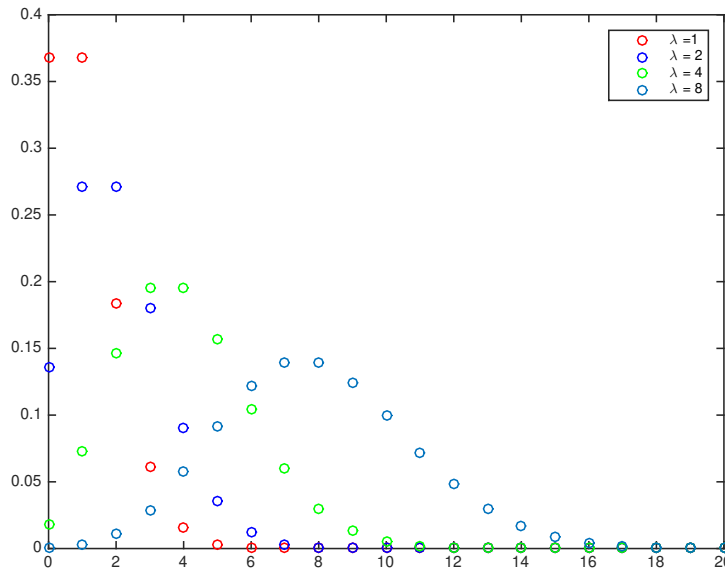
---

#### Example 1.1 (Photon Noise)

Light coming out of a laser pointer can be considered as a stream of photons. The detection or 'counting' of individual photons can be treated as independent random events that have a Poisson distribution. When  $N$  is the number of photons measured by a detector (a CCD camera or a photon multiplier tube) over a time interval  $t$  the Poisson distribution describing the number of counted photons is given as:

$$\Pr(N = k) = \frac{(\lambda t)^k e^{-\lambda t}}{k!} \quad (1.1)$$

Here  $\lambda$  represents the expected number of photons per unit time interval. In figure 1.1 we display the Poisson distribution (1.1) for different values of  $\lambda$ . For the



**Figure 1.1:** The Poisson distribution displaying the Probability of the number of photons  $k$  collected by a detector for different values  $\lambda$ , with  $\lambda$  the expected number of photon count per unit time interval.

Poisson distribution the variance of the number of photons over a time interval

$t$  is equal to,

$$\text{Var}[N] = \lambda t \quad (1.2)$$

This means that the 'strength' of the photon noise does depend on the length of the time interval used for counting photons. As the number of photons grows by  $\lambda t$ , the signal to noise ratio defined as the ratio between this number and the standard deviation equals

$$\text{SNR} = \sqrt{\lambda t} \quad (1.3)$$

This means that a longer exposure time will lead to a better signal to noise ratio. (Do you recognize this phenomena when taking a picture when it is dark and you do not have a flashlight?)

By an application of the central limit theorem to larger and larger photon counts, the Poisson distribution becomes Gaussian. This is why photon noise is also often modeled using a Gaussian distribution.

---

When observing a star during a clear sky at night we may observe light flickering. If we record this phenomena with a CCD camera, the photon count by a single pixel will change over time (even assuming a fixed exposure time). As illustrated by Example 1.1 the number of photons counted at each time instance is a sample of a random variable with a Poisson or a Gaussian distribution.

A second example of the relevance of stochastic processes for physics is Brownian motion.

---

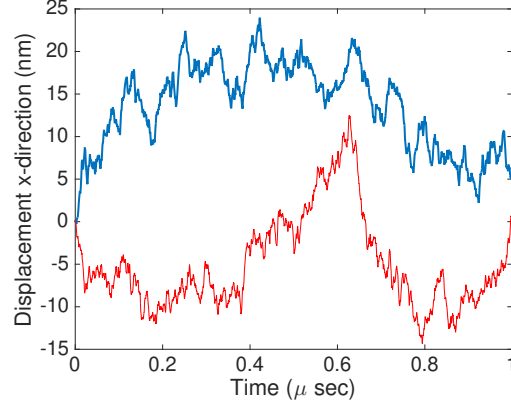
### Example 1.2 (Brownian Motion)

In modern physics Brownian motion is playing an increasingly important role. Particularly of interest is the Brownian motion in non-equilibrium systems related to the transport of molecules and cells in biological systems. Important examples [2] include Brownian motors, active Brownian motion of self-propelled particles, hot Brownian motion and Brownian motion in shear flows.

If we would observe and record the position of a particle in the  $x$ -direction during a certain time interval with a certain (constant) sampling frequency (e.g. 100 MHz), normalizing the position at time instant 0 as the zero position, then a sequence as displayed by the blue curve in Figure 1.2 would result. However, if we would repeat the experiment under the same experimental conditions (same temperature, etc.) a second time, a different graph will result. This is e.g. displayed by the red curve in Figure 1.2. In the simulation use has been made of a discretization of the Langevin equation discussed in Example 1.5 for the values of the quantities provided in Table 1.1.

---

Example 1.2 demonstrated that each observation of a time sequence over a time interval of the same length under identical experimental conditions of the position of a particle in the  $x$ -direction resulted in a different time record (or curve in the figure). As such it is concluded that the experiment is non-repeatable. This property of *non-repeatability* is characteristic of a stochastic process.



**Figure 1.2:** . Two time records of the  $x$ -displacement of a particle with mass  $m$  in a substance resulting in a friction coefficient  $\gamma$  at temperature  $T$  for  $1\mu s$  at a sampling rate of 100 MHz. Each time record the observation is done under identical experimental conditions with a normalization of the time axis and initial position at zero.

**Table 1.1:** The physical quantities and their values used in the simulation of the Langevin Equation (1.9).

| Quantity   | Meaning                 | Value                 | unit              |
|------------|-------------------------|-----------------------|-------------------|
| $N$        | number of samples       | $10^5$                | /                 |
| $\Delta t$ | sampling period         | $10^{-8}$             | s                 |
| $R$        | particle radius         | $10^{-6}$             | m                 |
| $k_B$      | Boltzmann constant      | $1.38 \cdot 10^{-23}$ | J/K               |
| $T$        | Temperature             | 300                   | K                 |
| $\eta$     | fluid viscosity (water) | $10^{-3}$             | Pa s              |
| $\rho$     | particle density        | $2.6 \cdot 10^3$      | kg/m <sup>3</sup> |
| $\gamma$   | friction coefficient    | $6\pi R\eta$          | Pa m s            |

### 1.3 In the footsteps of Carl Friedrich Gauss

The founding father of the field of stochastic signals and data driven modeling may be considered to be Carl Friedrich Gauss. He laid the foundations of data driven modeling with his famous least squares method. In the following Example 1.3 the least squares method is illustrated on a simple modeling problem. In the development of the least squares method Gauss developed a probability theory for the errors that describe the discrepancy between model and reality, such as indicated by the term  $e(t_i)$  in Example 1.3. Gauss introduced the normal distribution function as a natural way in which errors of observations occur. This choice of distribution is essential as it led to quadratic optimization problems that have an analytical solution. It enabled Gauss to find numerical solutions to simple parameter estimation problems by (tedious) hand calculations. This technique was an important asset in a time when digital computers were non-existent that was key to Gauss' success over many of his contemporaries.



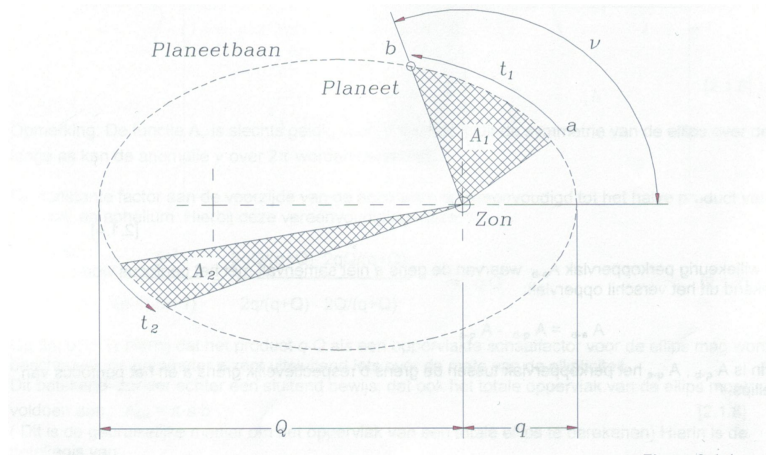
Carl Friedrich Gauss was a German mathematician and physicist born on April 30, 1777. He has had an immense influence on many fields of mathematics and physics. At the age of 23 he was challenged by the measurements of the Italian Monk Giuseppe Piazzi made about the orbit of the planet Ceres. Piazzi was only able to make observations for 41 days early in 1801. Then the planet disappeared, challenging many scientists of that day to predict the reappearance of Ceres. Gauss used the method of least squares to determine features of the planet's orbit. Based on these calculations he was able to predict the reappearance of Ceres to within half a degree.

Though Gauss claimed to have used the method of least squares for the accurate prediction of the orbit of Ceres, the principles of (linear) least squares were first published by the French mathematician Legendre in 1805. Gauss claimed that he had discovered the method earlier than 1801 but did not publish it because he did not consider it very important. Gauss then later described the least squares method in the papers *Theoria combinationis observationum erroribus minimis obnoxiae I and II* of 1821 and 1823.

---

**Example 1.3 (Predicting the re-appearance of a planet)**

A number of articles have appeared trying to discover the original calculations of Gauss in the prediction of the orbit of Ceres [3]. In this example, we treat a simplified scenario in predicting orbits of planets to illustrate the method of least squares in a deterministic context. The illustration is based on [4].

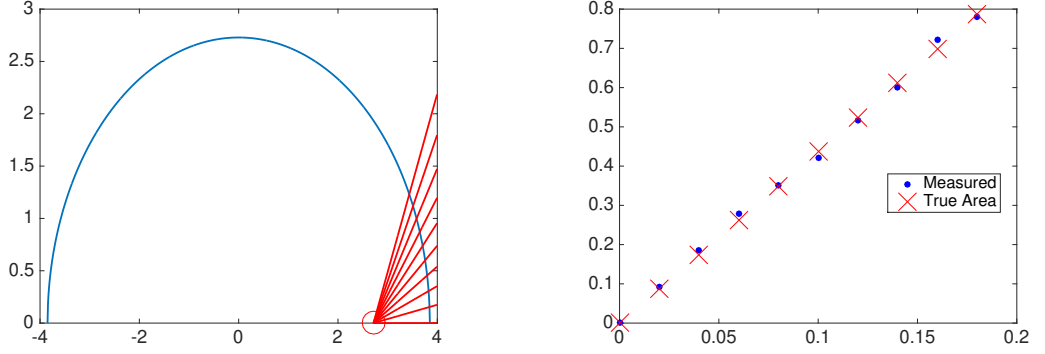


**Figure 1.3:** A two dimensional plot of the elliptical orbit of a planet with the sun as one of its focal points. The two indicated sizes  $q, Q$  and the eccentricity  $e$  enable to uniquely define the ellipse.

Consider Figure 1.3 and assume that from the comet Shoemaker-Levy the following measurements of the area  $A_m(t)$  covered of the planet and the corresponding time instances  $t$  are as displayed in Tabel 1.2. The covered area is pictured in Figure 1.4.

| Data Point Number | Area covered $A_m(t_i)$<br>(Normalized Unit) | Elapsed time $t_i$<br>(Normalized Unit) |
|-------------------|--|---|
| 1                 | 0.0075                                       | 0                                       |
| 2                 | 0.0855                                       | 0.02                                    |
| 3                 | 0.1837                                       | 0.04                                    |
| 4                 | 0.2545                                       | 0.06                                    |
| 5                 | 0.3355                                       | 0.08                                    |
| 6                 | 0.4227                                       | 0.10                                    |
| 7                 | 0.5292                                       | 0.12                                    |
| 8                 | 0.6100                                       | 0.14                                    |
| 9                 | 0.6972                                       | 0.16                                    |
| 10                | 0.8007                                       | 0.18                                    |

**Table 1.2:** The numerical data about the covered area and elapsed time of the Comet Shoemaker-Levy [4].



**Figure 1.4:** (Left) Pictorial representation of the 10 ‘virtual’ observations of the position of the comet Shoemaker-Levy. The red small circle represents the position of the sun and the crossing of the red lines and the blue orbit represent the respective positions of the comet with the starting point in the vicinity of the point (4,0). The planet moves counter-clockwise. (Right) The measured covered area versus time: with  $\times$  — the true area covered and with  $\bullet$  — the measured area covered. The small discrepancies between red crosses and blue dots is due to additive (measurement) errors.

Following Kepler’s second law, the area covered until time  $t$  and denoted by  $A(t)$  is proportional to the elapsed time  $t$  as given by the following model:

$$A(t_i) = kt_i \quad (1.4)$$

Here  $k$  is an unknown constant, characteristic for the comet under study. However when we derive a value of  $k$  at each time instant  $t_i$  for which a measurement  $A_m(t_i)$  is available from the model (1.4), each time instance we would observe a different  $k$ . This problem of non-uniqueness is in this example a direct consequence of the measurement and calculating errors causing the equality (law) in (1.4) to be violated. In order to resolve this violation, the brilliant idea of Gauss was to introduce an error term  $e(t)$  in (1.4) as follows:

$$A_m(t_i) = kt_i + e(t_i) \quad i = 1 : 10 \quad (1.5)$$

The goal is now to determine an ‘averaged’ parameter  $k$  such that the ‘sum of squared errors’  $\sum_{i=1}^N e(t_i)^2$  (for  $N = 10$ ) is minimized. This is a deterministic optimization problem denoted as follows:

$$\min_k \sum_{i=1}^{10} \left( A_m(t_i) - kt_i \right)^2 \quad (1.6)$$

The necessary (and sufficient) conditions that characterize the optimum  $\hat{k}$  are:

$$\left. \frac{\partial \sum_{i=1}^{10} \left( A_m(t_i) - kt_i \right)^2}{\partial k} \right|_{k=\hat{k}} = 0 \Rightarrow \sum_{i=1}^{10} A_m(t_i) t_i - \hat{k} \sum_{i=1}^{10} t_i^2 = 0 \quad (1.7)$$

For the given data in Table 1.2 this yields  $\hat{k} = 4.3687$ . This combined with the total area covered of 32.9985 yields a period of the orbit of 7.5535 while the true period is 7.5519.

---

## 1.4 Objectives of the Course

This course provides an introduction to stochastic processes for physics students. The course has 4 objectives. These are discussed in the following subsections.

### 1.4.1 Objective 1: Description of stochastic processes

Example 1.2 illustrates that stochastic processes are non-repeatable. This is contrary to deterministic signals that may be reproduced exactly, e.g. given by an exact mathematical prescription.

---

#### Example 1.4 (Deterministic signal)

One example of a discrete harmonic sequence is:

$$x(n) = A \sin(n\omega_0 + \phi) \quad (1.8)$$

with  $A, \omega_0$  fixed positive real numbers and  $\phi$  a fixed real number that satisfies  $-\pi \leq \phi < \pi$ . The prescription given by the mathematical formula (1.8) determines the signal  $x(n)$  exactly for each (integer) value of the sample index  $n$ .

---

A statistical framework will be developed to describe stochastic processes via statistical quantities. These describing quantities are or may be deterministic. An example of such a statistical quantity is the mean or first order statistical moment of the stochastic processes as a function of time. The main statistical quantity used to describe stochastic processes in this course are second order statistical moments such as the Auto- and Cross correlation functions. After the definition of the statistical characteristics that will be used in this course to characterize stochastic processes in Chapter 5 we will also discuss in this chapter how and under which conditions some of these statistical characteristics can be retrieved from a single time record of the stochastic process.

### 1.4.2 Objective 2: Generation of stochastic processes

Stochastic processes, such as the Brownian motion in Example 1.2, may be modeled by a stochastic differential equation known as the Langevin equation. This is further illustrated in Example 1.5.

---

#### Example 1.5 (Langevin)

The 1-dimensional dynamics (in the  $x$ -direction) of a spherical Brownian particle with mass  $m$  in a gas can be described by the Langevin equation [2, 5]:

$$m \frac{d^2 x(t)}{dt^2} + \gamma \frac{dx(t)}{dt} = \sqrt{2k_B T \gamma} w(t) \quad (1.9)$$

The variables in this equation are listed in Table 1.3. The driving force  $w(t)$  of this Langevin equation is so-called continuous white noise and turns the equation

**Table 1.3:** The quantities in the Langevin equation (1.9).

| Quantity    | Meaning   |
|-------------|---|
| $m$         | mass of the particle  |
| $\gamma$    | friction coefficient  |
| $x(t)$      | displacement of the particle in the $x$ -direction at time instant $t$                      |
| $k_B$       | Boltzmann constant  |
| $T$         | Temperature   |
| $w(t)$      | normalized white noise: $\forall t, t' : E[w(t)] = 0 \quad E[w(t)w(t')] = \delta(t - t')$ . |
| $\delta(t)$ | the Dirac delta function  |

(1.9) into a stochastic differential equation. The solution of such an equation is generally done by discretization as discussed in [6]. Here also a discrete variant of the white noise signal is discussed. This discretization results in the following difference equation:

$$x_d(n) + a_1 x_d(n-1) + a_2 x_d(n-2) = b_0 w_d(n) \quad (1.10)$$

where  $x_d(n)$  is the approximation due to discretization and definition of the discrete white noise signal  $w_d(n)$  of the  $x$ -displacement at time instance  $n\Delta T$ . For simplicity we normalize the time scale, allowing us to consider the sample interval  $\Delta T$  to be 1. For a definition of a discrete white noise signal we refer to Chapter 3.

For that reason we study in Chapter 6 how the statistical characteristics of a stochastic process is changed when filtering that stochastic process with a Linear Time Invariant (LTI) filter.

Inversely we are interested in finding an LTI filter and its input such that the filtered output has a statistical characteristic that matches a given characteristic. These problems both in the time and frequency domain are discussed in Chapter 7.

### 1.4.3 Objective 3: Parameter Estimation

A more fundamental question is how to find models such as the Langevin equation (1.9) from just one time recording of the displacement over a finite time interval. This fundamental problem also belongs to the class of so-called *inverse problems*. A simplified analysis of this fundamental question is performed in Chapter 8, where we develop the linear least squares and outline its use in addressing simple inverse problems. The outline is a generalization of the simple Example 1.3, now treating the signals and the criterion to be optimized in a stochastic setting.

### 1.4.4 Objective 4: Optimal Filtering

A problem of general interest is to estimate one signal that cannot be measured by filtering another signal that can be measured assuming that both signals bear

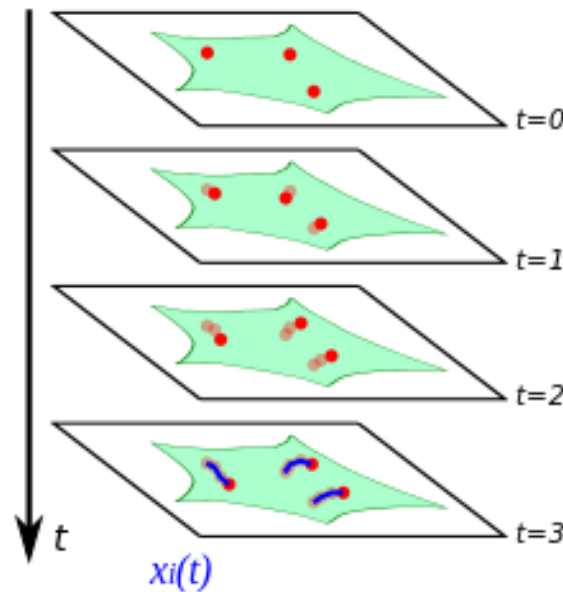


some (statistical) relationship. The determination of the filter such that a chosen cost function is optimized is indicated by *optimal filtering*. Four different optimal filtering problems, to be briefly discussed in Section 1.5 will be considered in this course. A physical example of particle tracking is outlined next.

---

#### Example 1.6 (Single Particle Tracking)

Tracking of particles is of interest e.g. to understand the cellular kinetics of proteins like HIV-1. This is used for example in studying the efficacy of drugs to combat certain diseases. A classical approach in particle tracking aims at localizing the same particle in a sequence of frames, as illustrated in Figure 1.5. However, novel approaches make use of parametric prediction methods [7]. If we denote the position of a particle at time instance  $t$  in the  $x$ -direction by  $x(t)$ , then a challenge with novel parametric prediction methods (within the scope of this course) is to find an LTI filter that filters  $x(t)$  such that the output of the filter is a prediction of the position at the next time instant  $t + \Delta T$ . Using this prediction based on information up to time instant  $t$ , will enable a better localization of the particle the next time a frame is taken at time instant  $t + \Delta T$ .



**Figure 1.5:** Single Particle Tracking by localizing individual particles in a sequence of images. The first frame on top contains 3 particles (indicated by the red dots), the second frame below contains the new location (yellow dots) plus the old location of the three particles, and the bottom frame glues these estimated location of each particle together in a 'blue trajectory'.

Picture from: [http://en.wikipedia.org/wiki/Single\\_particle\\_tracking](http://en.wikipedia.org/wiki/Single_particle_tracking)

---

---

## 1.5 Applications of Optimal Filtering

Four problems are addressed in this course that are illustrations of the general filtering problem stated in Objective 4 of this course to retrieve (an unmeasurable stochastic) signal by filtering another measurable signal when both signals are (statistically) related. These four problems are briefly postulated next.

### 1.5.1 Denoising

Sensors used to record signals have finite precision. In addition they capture unwanted phenomena that we generally call noise. Consider e.g. the ECG (electrocardiogram) signal. A typical recording of an ECG is shown in figure 1.6 on the left. Using digital hardware such recording consists of an ordered (in time) sequence of (amplitude) samples. Such signals will be called *discrete (time) signals*. The electrodes used to capture this signal introduce noise coming from various sources. Such as due to poor electrode-skin contact, interference with other sources, electrical noise, etc. In order to denote this mathematically, let the recorded ECG signal be denoted by  $x(n)$  and let the noise free ECG signal be denoted by  $d(n)$ . Then under the assumption that the noise  $v(n)$  is *additive* the following signal relationship holds:

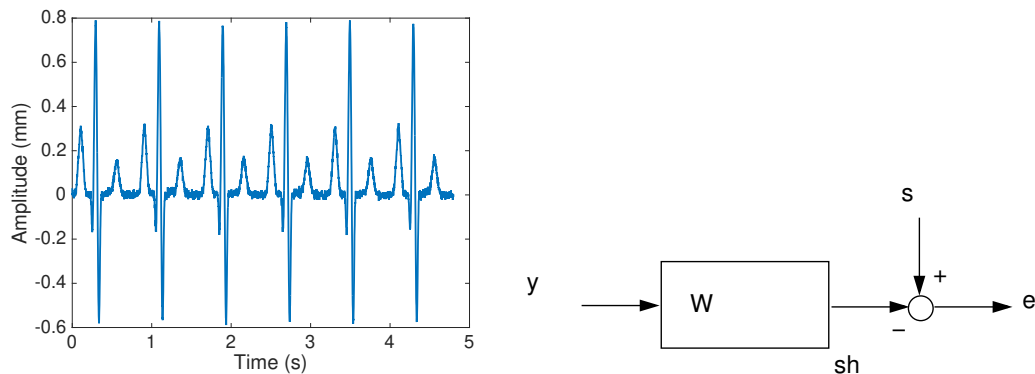
$$x(n) = d(n) + v(n) \quad (1.11)$$

In order to retrieve an estimate of the noise-free ECG signal  $d(n)$  one can aim to design a Linear Time-invariant (LTI) filter  $W(z)$  to filter the measured signal  $x(n)$  as depicted in figure 1.6 on the right. The filtered signal is denoted here by  $\hat{d}(n)$  and the goal that we will develop in this course is to determine the parameters of the filter  $W(z)$  such that the error signal  $e(n)$  is *small*.

When the signals under investigation are stochastic processes a more precise formulation will be given in this course in Chapter 9 on what is indicated by *small*. In addition to this formulation of a criterion to express what is statistically *small*, the solution of this problem motivates the need to (a) describe statistical relationship between two stochastic processes via e.g. their cross-correlation function, (b) describe how these statistical characteristics change by LTI filtering and (c) develop a methodology to optimize the statistical criterion and analyse the quality of the optimum.

### 1.5.2 Deconvolution

A generalization of the denoising problem is when the signal  $x(n)$  is a filtered (by a known LTI filter) version of the signal  $d(n)$  plus noise  $v(n)$ . Such problems occur for example in digital image processing where the filtering corresponds to a blurring operation on the original image. The filtering is now not a temporal but a spatial filtering and the input signal is a (digitized) image. An example is given next.



**Figure 1.6:** Typical ECG recording disturbed by additive noise.

---

### Example 1.7 (Image Blurring)

Consider the TUD logo depicted in Figure 1.7 on the left. When take a picture with a camera with a linear motion of 11 pixels in the horizontal direction, a blurred picture results as depicted in Figure 1.7 on the right. This result is obtained by filtering the original image with a 2D filter with impulse response generated by the Matlab command `fspecial('motion', 11, 0)`.



**Figure 1.7:** The original TUD logo (on the left) and the blurred TUD logo taken by a camera undergoing a linear horizontal motion of 11 pixels.

---

In this course we restrict to 1D filters only and generally this one dimension is time, but it may also be a spatial dimension.

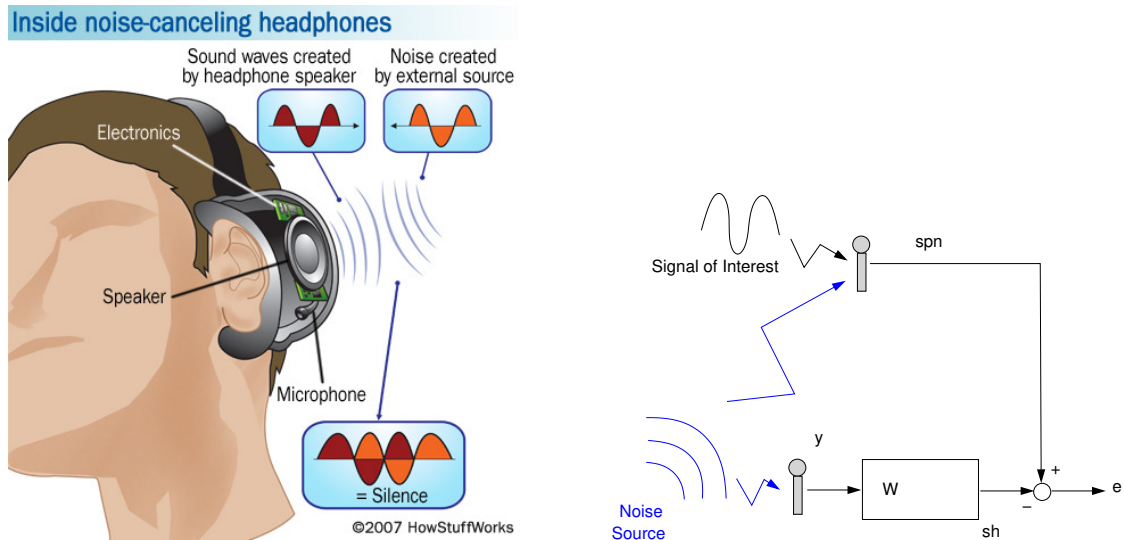
### 1.5.3 Prediction

The prediction problem has been illustrated by the Example 1.6.

### 1.5.4 Active Noise Cancellation

Active noise control is trying to control the acoustics in a closed environment. An artist's view of active noise control in a headset is depicted in the left hand of Figure 1.8. Here the goal is to activate the speaker in the headset such that it cancels as much as possible the unwanted sound  $v_1(n)$  picked up by the microphone through which the person in the noisy environment is communicating

(the speaker micro). In this course we will look at the simple problem of designing a filter  $W(z)$  in the block-scheme in the right hand of Figure 1.8. The output of the filter  $W(z)$  should represent the signal  $\hat{v}_1(n)$  that is reproduced by the speaker in the headset to cancel  $v_1(n)$  (as much as possible). The filter  $W(z)$  uses as input the signal  $v_2(n)$  that is assumed to ‘well’ related to the unwanted noise. This may be achieved by a second micro in the headset that located far enough from the speaker micro.



**Figure 1.8:** (Left) Active Noise cancellation in a headset. The components are two microphones: one near the mouth of the speaker picking up the sound of the speaker and the noisy environment, indicated by the signal  $d(n) + v_1(n)$  in the right figure, and a second in the headset picking up mainly sound from the environment, indicated by the signal  $v_2(n)$  in the right figure. The headset also has in addition to electronics a speaker that produces the sound  $\hat{v}_1(n)$ . (Right) Block Schematic representation of the Left.

The left figure is taken from <http://electronics.howstuffworks.com/gadgets/audio-music/noise-canceling-headphone3.htm>.

To analyse and address this problem it is again necessary to express the notion ‘as much as possible’ as a mathematical criterion. This criterion depends on the observed signals and on the parameters of the filter  $W(z)$  to be tuned. Also here the problems (a-c) mentioned in the last paragraph of the denoising problem are relevant.

## References

- [1] R. Brown, “A brief account of microscopical observations made in the months of june, july and august, 1827, on the particles contained in the pollen of plants; and on the general existence of active molecules in organic and inorganic bodies,” *Philosophical Magazine*, vol. 4, pp. 161–173, 1828.

- [2] T. Li and M. Raizen, "Brownian motion at short time scales," arXiv:1211.148v1.
- [3] J. Tennenbaum and B. Director, "How gauss determined the orbit of ceres," *Fidelio*, Summer 1998.
- [4] R. Ott, "Keplers orbits onconventional: Mathematical treatment for an alternative derivation of keplers second law," 1998. <http://www.dutch.nl/rcott/astrom.htm>.
- [5] P. Langevin, "Sur la théorie du mouvement brownien," *C. R. Acad. Sci. Paris*, vol. 146, pp. 530–533, 1908.
- [6] G. Volpe and G. Volpe, "Simulation of a brownian particle in an optical trap," *American J. Phys*, vol. 81, no. 3, pp. 224–230, 2013.
- [7] I. Smal, I. Grigoriev, A. Akhmanova, W. J. Niessen, and E. Meijering, "Accurate estimation of microtubule dynamics using kymographs and variable-rate particle filters," in *31st Annual International Conference of the IEEE EMBS*, pp. 1012–1015, 2009.

---

## Exercises

- Exercise 1.1** Provide one additional example to those given in Section 1.2 of physical problem where stochastic signals are present.
- Exercise 1.2** Give an additional application to the application of Optimal filtering to those given in Section 1.5.



# Chapter 2

## Signals and Systems

**After studying this chapter you can:**

- describe discrete-time signals with the use of unit pulses.
- recall the definition of Discrete-Time Fourier Transform (DTFT) and understand that it is a special case of the z-transform.
- use the properties of the DTFT and the z-transform as listed in Tables 2.2 and 2.3
- derive the Region of convergence (ROC) of rational transfer functions.
- reuse the notions about Linear-Time Invariant systems, stability, causality and minimum-phase.
- perform a splitting in causal and strictly anti-causal part of a mixed causal and anti-causal rational transfer function via partial fraction expansion.

---

## 2.1 Introduction

In this chapter we first review some basic concepts of discrete-time signals both in the time domain and in the frequency domain. As signals interact with each other via systems a review is given secondly about notions to define and analyse discrete-time systems. The review is concise and meant to introduce the terminology and notation used throughout these course notes. For a more elaborate exposure on signals and systems we refer to books, like [1, 2].

The organization of the chapter is as follows. Section 2.2 contains a review of the description of discrete time signals in time-, Fourier domain and with the help of the z-transform. Relevant concepts, such as stability, causality and minimum-phase property of dynamical Linear Time Invariant (LTI) systems are reviewed in Section 2.3. The final section 2.4 gives an example of a simple inverse problem.

---

## 2.2 Discrete-Time Signals

### 2.2.1 Definition and examples

A discrete-time signal is an indexed sequence of numbers or vectors. The numbers or the entries of the vectors are either real or complex. Though we mainly focus on time-domain signals, the indexing may be with respect to time or space. Spatial indexing typically occurs when representing a digital image in terms of its pixels and their coordinates in 2D. We restrict to one-dimensional signals where the indexing is represented by integer values.

A common way discrete-time signals occur is by sampling continuous-time signals. Let the continuous-time signal be denoted by  $x_c(t)$  and let the sampling be done at the equidistant time instances  $t = n\Delta T$ , then the relationship between the discrete-time and continuous time signals can be denoted as:

$$x(n) = x_c(n\Delta T) \quad n \in \mathbb{Z} \quad (2.1)$$

Special discrete-time signals are defined in Table 2.1.

| Signal                          | Mathematical Notation | definition  |
|---------------------------------|-----------------------|---|
| unit pulse (or Kronecker delta) | $\delta(n)$           | $\delta(n) = \begin{cases} 1 & \text{for } n = 0 \\ 0 & \text{for } n \neq 0 \end{cases}$ |
| unit step                       | $u(n)$                | $u(n) = \begin{cases} 1 & \text{for } n \geq 0 \\ 0 & \text{for } n < 0 \end{cases}$      |
| complex exponential             | $e^{jn\omega}$        | $e^{jn\omega} = \cos(n\omega) + j \sin(n\omega)$  |

**Table 2.1:** Three specific discrete-time signals.

The signals defined in Table 2.1 have infinite length. They are also referred to as two-sided sequences with the index  $n = 0$  at the middle of the sequence.



When a sequence  $x(n)$  would be defined only for positive indices  $n \geq 0$  it is referred to as a right sided sequence and when defined for negative indices only a left sided sequence. The two-sided, left- and right-sided sequences are all infinite sequences.

Discrete time-signals can also be defined over a finite interval. In that case they are indicated as *finite sequences*.

The unit pulse signal  $\delta(n)$  (see Table 2.1) can be used in the definition of a discrete-time sequence. Consider  $x(n)$  to be an infinite sequence then it can be represented as,

$$x(n) = \sum_{k=-\infty}^{\infty} x(k)\delta(n-k) \quad (2.2)$$

In this way we can embed a finite sequence into an infinite one as follows. Consider the finite sequence  $x_N(n)$  defined for  $n \in [1, N]$  then the infinite sequence  $x(n)$  embedding  $x_N(n)$  can be defined as:

$$x(n) = \sum_{k=-\infty}^{\infty} x_N(k)\delta(n-k) \quad \text{with } x_N(k) = 0 \quad k \notin [1, N] \quad (2.3)$$

## 2.2.2 The Discrete-Time Fourier Transform

Discrete-time signals can be analysed in the frequency domain by the use of the Discrete-Time Fourier Transform (DTFT). When the discrete time signal is assumed to be a sampled continuous-time signal defined at the time instances  $n\Delta T$  as  $x(n\Delta T)$  and when the resulting sequence is *absolute summable*, i.e.

$$\sum_{n=-\infty}^{\infty} |x(n\Delta T)| < \infty \quad (2.4)$$

then the DTFT of the signal  $x(n\Delta T)$  is defined as:

$$X(e^{j\omega}) = \sum_{n=-\infty}^{\infty} x(n\Delta T)e^{-j\omega n\Delta T} \quad \omega \in \mathbb{R}$$

In the major part of these course notes we will not consider the connection with sampled continuous-time signals and hence will consider  $\Delta T = 1$ . In that case the DTFT is given as in (2.5). Here we also introduce the notation used to indicate this transform.

$$X(e^{j\omega}) = \mathcal{F}(x(n)) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n} \quad \omega \in \mathbb{R} \quad (2.5)$$

This definition shows that the DTFT transforms a discrete sequence into a continuous function, i.e. the function  $X(e^{j\omega})$  is continuous in its argument  $\omega$ . Since  $e^{-j\omega n}$  is *periodic* with a period  $2\pi$ , the DTFT is also periodic with period  $2\pi$ . In general the DTFT is a complex-valued function of  $\omega$  and therefore  $X(e^{j\omega})$  can be represented in terms of its magnitude and phase,

$$X(e^{j\omega}) = |X(e^{j\omega})|e^{j\angle(X(e^{j\omega}))} \quad (2.6)$$

The absolute summable condition (2.4) is a sufficient condition for the existence of the DTFT. However the use of generalized functions allows to express the DTFT of signals that violate the condition (2.4). When  $\delta(\omega)$  denotes an impulse [3] at frequency  $\omega = 0$  then the DTFT of a complex exponential is given as:

$$\mathcal{F}(e^{jn\omega_0}) = 2\pi\delta(\omega - \omega_0) \quad |\omega| < \pi \quad (2.7)$$

In the same way the DTFT of the non-absolute summable unit step function  $u(n)$ , defined in Table 2.1, can be expressed with the use of the generalized impulse function [2] as,

$$\mathcal{F}(u(n)) = \frac{1}{1 - e^{-j\omega}} + \pi\delta(\omega) \quad |\omega| < \pi \quad (2.8)$$

For the case  $x(n) \in \mathbb{R}$ , then  $X(e^{j\omega})$  is *conjugate symmetric*, i.e.

$$X(e^{j\omega}) = X^*(e^{-j\omega}) \quad (2.9)$$

The DTFT is an invertible transformation. This means that the original signal  $x(n)$  may be calculated when its Discrete Fourier Transform  $X(e^{j\omega})$  is given. This is provided by the inverse DTFT as follows,

$$x(n) = \mathcal{F}^{-1}(X(e^{j\omega})) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega}) e^{jn\omega} d\omega \quad (2.10)$$

Finally with the definition of the convolution between 2 infinite sequences  $x(n)$  and  $y(n)$  given as,

$$x(n) \star y(n) = \sum_{k=-\infty}^{\infty} x(k)y(n-k) \quad (2.11)$$

we summarize in Table 2.2 a number of properties of the DTFT.

| Property       | Time signal                            | DTFT   |
|----------------|--|--|
| Original       | $x(n)$<br>$y(n)$                       | $X(e^{j\omega}) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n}$<br>$Y(e^{j\omega}) = \sum_{n=-\infty}^{\infty} y(n)e^{-j\omega n}$ |
| Linearity      | $ax(n) + by(n)$                        | $aX(e^{j\omega}) + bY(e^{j\omega})$ for $a, b \in \mathbb{C}$  |
| Time Shift     | $x(n - \alpha)$                        | $e^{-j\omega\alpha} X(e^{j\omega})$ for $\alpha \in \mathbb{Z}$  |
| Conjugation    | $x^*(n)$                               | $X^*(e^{-j\omega})$  |
| Time Reversal  | $x(-n)$                                | $X(e^{-j\omega})$  |
| Convolution    | $x(n) \star y(n)$                      | $X(e^{j\omega})Y(e^{j\omega})$   |
| Multiplication | $x(n)y(n)$                             | $\frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\beta})Y(e^{j(\omega-\beta)})d\beta$   |
| Parseval       | $\sum_{n=-\infty}^{\infty} x(n)y^*(n)$ | $\frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega})Y^*(e^{j\omega})d\omega$   |

**Table 2.2:** Some Properties of the DTFT.

### 2.2.3 The z-Transform

A generalization of the DTFT is the z-transform. The z-transform of a signal  $x(n)$  and its notation are defined as,

$$X(z) = \mathcal{Z}(x(n)) = \sum_{n=-\infty}^{\infty} x(n)z^{-n} \quad z = re^{j\omega} \in \mathbb{C} \quad r, \omega \in \mathbb{R} \quad (2.12)$$

This definition shows that the DTFT is a special case of the z-transform. Namely that when  $z = e^{j\omega}$  the z-transform becomes the DTFT. The z-transform of signals may exist for which the DTFT does not.

The existence of the z-transform is associated with the values of  $z \in \mathbb{C}$  for which the following infinite sum

$$\sum_{n=-\infty}^{\infty} |x(n)||z|^{-n} \quad (2.13)$$

is *bounded*. The values of  $z \in \mathbb{C}$  for which the sum in (2.13) converges is called the *region of convergence (ROC)* of the z-transform.

We now discuss three special sequences and their ROC:

1. *Finite length sequences*: When  $x_N(n)$  is a finite time sequence defined for  $n \in [0, N-1]$  its z-transform is a finite series in  $z$  or a polynomial in  $z$ . The ROC of this z-transform is the whole complex plane except possibly  $z = 0$ .
2. When  $x(n)$  is a right-sided sequence only defined for  $n \geq 0$  and if the circle  $|z| = R_-$  is in the ROC, then all finite values of  $z$  for which  $|z| > R_-$  will also be in the ROC:

$$|z| \geq R_- > 0$$

3. When  $x(n)$  is a left-sided sequence only defined for  $n \leq 0$  and if the circle  $|z| = R_+$  is in the ROC, then all values of  $z$  for which  $0 < |z| < R_+$  will also be in the ROC:

$$|z| \leq R_+, \quad R_+ > 0$$

For a general two-sided sequence that has the circle  $|z| = R_0$  in the ROC, the ROC of its z-transform is a region in the complex plane  $\mathbb{C}$  contained in between 2 circles that includes the circle  $|z| = R_0$ . When this region contains the unit circle in  $\mathbb{C}$ , the DTFT of the two-sided sequence exists.

As the DTFT, the z-transform has a number of related properties, summarized in Table 2.3.

**Remark 2.1** (Inverse z-transform). *By the use of the so-called double-sided z-transform (2.12) (with the lower bound and upper bound of the summation index going to resp.  $-\infty$  and  $\infty$ ), the inverse of the z-transforms that are given in the third column of Table 2.3 yields the corresponding expressions in the second column.*

The z-transform of four special signals is given in Table 2.4.

A final result that will be used in Chapter 9 is a specific case of Cauchy's Integral formula applied to the z-transform in (2.12) [7]. Here we will make

| Property      | Time signal       | z-transform  |
|---------------|-------------------|--|
| Original      | $x(n), y(n)$      | $X(z) = \sum_{n=-\infty}^{\infty} x(n)z^{-n}, Y(z) = \sum_{n=-\infty}^{\infty} y(n)z^{-n}$ |
| Linearity     | $ax(n) + by(n)$   | $aX(z) + bY(z)$ for $a, b \in \mathbb{C}$  |
| Time Shift    | $x(n - \alpha)$   | $z^{-\alpha} X(z)$ for $\alpha \in \mathbb{Z}$   |
| Conjugation   | $x^*(n)$          | $X^*(z^*)$   |
| Time Reversal | $x(-n)$           | $X(z^{-1})$  |
| Convolution   | $x(n) \star y(n)$ | $X(z)Y(z)$   |

**Table 2.3:** Some Properties of the z-transform.

| Signal                                | z-transform                       | ROC                     |
|---------------------------------------|-----------------------------------|-------------------------|
| $\Delta(n)$                           | 1                                 | $\mathbb{C}$            |
| $a^n u(n) \quad a \in \mathbb{R}$     | $\frac{1}{1-az^{-1}}$             | $ z  > a$               |
| $-a^n u(-n-1) \quad a \in \mathbb{R}$ | $\frac{1}{1-az^{-1}}$             | $ z  < a$               |
| $a^{ n }$                             | $\frac{1-a^2}{(1-az^{-1})(1-az)}$ | $a <  z  < \frac{1}{a}$ |

**Table 2.4:** The z-transform and its ROC of four special signals.

use of the symbol  $\oint_{\Gamma}$  to indicate the circle integral over a closed-curve  $\Gamma$  in the complex plane. This closed-curve will always be the unit circle in the complex plane in these course notes.

**Theorem 2.2** (Cauchy's Integral Formula (Simplified)). *Let  $X(z)$  be the z-transform as given in (2.12) with ROC containing the unit circle, denoted by the curve  $\Gamma$  in the complex plane, then,*

$$x(0) = \frac{1}{2\pi j} \oint_{\Gamma} \frac{X(z)}{z} dz$$

## 2.3 Discrete-Time Systems

### 2.3.1 Definition

A discrete-time system is a transformation of its input signal, denoted by  $x(n)$  in Figure 2.1, into its output signal, denoted by  $y(n)$  in Figure 2.1. The transformation is represented as,

$$y(n) = H[x(n)] \quad (2.14)$$



**Figure 2.1:** Schematic Representation of a Discrete-Time system represented by the operator or mapping  $H[-]$ , transforming its input  $x(n)$  to its output  $y(n)$ .

An example of a discrete-time system is given in Example 1.5 by the difference equation (1.10).

### 2.3.2 Linear Time-invariant (LTI) Discrete Time Systems

In this course we restrict to the special class of Linear, Time-invariant (LTI) systems.

**Definition 2.3** (Linearity). *A discrete-time system given by the operator or mapping  $H[-]$  is linear if,*

$$H[ax_1(n) + bx_2(n)] = aH[x_1(n)] + bH[x_2(n)] \quad a, b \in \mathbb{C} \quad (2.15)$$

**Definition 2.4** (Time-invariance). *A discrete-time system given by the operator or mapping  $H[-]$  that has the following input-output relation  $y(n) = H[x(n)]$ , is time-invariant if,*

$$H[x(n - \alpha)] = y(n - \alpha) \quad \alpha \in \mathbb{Z} \quad (2.16)$$

For an LTI system the response to any input sequence may be expressed as a convolution between this input sequence and the response to a unit pulse. The latter response is called the *unit pulse response* or *impulse response* of the LTI system. The convolution result can be deduced as follows. Let the LTI system be given by the operator or mapping  $H[-]$  have unit pulse response denoted as  $h(n)$  and consider the input  $x(n)$  represented as in (2.2) then by the Linearity property we can represent the output  $y(n)$  as,

$$y(n) = \sum_{k=-\infty}^{\infty} x(k)H[\Delta(n - k)]$$

By the time-invariance property the convolution result follows,

$$y(n) = \sum_{k=-\infty}^{\infty} x(k)h(n - k) = x(n) \star h(n) \quad (2.17)$$

By the convolution property of the DTFT and the z-transform we can formulate the Fourier transform, resp. the z-transform of the output, as a product. To illustrate this for the DTFT, let  $Y(e^{j\omega})$ ,  $X(e^{j\omega})$ ,  $H(e^{j\omega})$  resp. denote the DTFT transforms of the signals  $y(n)$ ,  $x(n)$  and the unit pulse response  $h(n)$ . Here we assume that these transforms exist. Then the convolution property of the DTFT in Table 2.2 shows that the DTFT of the convolution in (2.17) becomes,

$$Y(e^{j\omega}) = H(e^{j\omega})X(e^{j\omega}) \quad (2.18)$$

Likewise using the convolution property of the z-transform in Table 2.3 we have the z-transform of the convolution in (2.17) becomes,

$$Y(z) = H(z)X(z) \quad (2.19)$$

The z-transform  $H(z)$  of an LTI system with impulse response  $h(n)$  is called the *transfer function* of the LTI system.

The representation of LTI systems via their impulse response shows that LTI systems can be interpreted as a signal. Further the product relationships (2.18) or (2.19) show that the 'role' of the input and impulse response can be interchanged. For example in (2.19) the output can be interpreted as well as the convolution of the 'input signal'  $h(n)$  with a system with impulse response  $x(n)$ .

### 2.3.3 Stability, Causality and Minimum-phase LTI systems

For LTI systems there are 3 additional properties of particular interest for the rest of this course. These are:

1. Bounded-input, Bounded-output (BIBO) stability.
2. Causality and anti-causality
3. Inverse of a system and minimum-phase systems

These properties are discussed next.

**Definition 2.5** (BIBO stability). *A system given by the operator, or mapping  $H[-]$  is BIBO (Bounded-input, Bounded-output) stable if any bounded input  $u(n)$ , satisfying,*

$$\max_n |u(n)| < \infty$$

*yields a bounded output  $y(n) = H[u(n)]$ .*

A necessary and sufficient condition for LTI systems is given in the following Lemma.

**Lemma 2.6** (LTI BIBO stability). *An LTI system given by its impulse response  $h(n)$  is BIBO stable if and only if  $h(n)$  is absolute summable.*

*Proof.* The sufficiency is shown that under the condition of absolute summability any bounded input, yields a bounded output. For that purpose, consider the convolution (2.17) expressed as in Exercise 2.5, then the absolute value of  $y(n)$  satisfies,

$$|y(n)| = \left| \sum_{k=-\infty}^{\infty} h(k)x(n-k) \right| \leq \sum_{k=-\infty}^{\infty} |h(k)||x(n-k)|$$

When the input is bounded, there exists a real number  $M$  such that  $\max_n |x(n)| < M$ . Therefore, it follows from the above relationship that,

$$|y(n)| \leq M \sum_{k=-\infty}^{\infty} |h(k)|$$

and since  $h(n)$  was assumed to absolute summable, the output is bounded.

To proof the necessity condition we will show that there exists a bounded input, that yields an unbounded output when  $h(n)$  is not absolutely summable. Now consider the bounded input (for a fixed  $n$ ):

$$x(n-k) = \begin{cases} 1 & \text{for } h(k) \geq 0 \\ -1 & \text{for } h(k) < 0 \end{cases}$$

In that case the convolution (2.31) shows that,

$$|y(n)| = \left| \sum_{k=-\infty}^{\infty} h(k)x(n-k) \right| = \sum_{k=-\infty}^{\infty} |h(k)|$$

And the right hand side is not bounded. □

When the impulse response is absolute summable, it means that the unit circle is part of the ROC of the z-transform of the impulse response.

**Definition 2.7** (LTI Causal or Anti-Causal systems). *An LTI system given by its impulse response  $h(n)$  is (strict) causal if and only if  $h(n)$  is zero for  $n < 0$  ( $n \leq 0$ ).*

*An LTI system given by its impulse response  $h(n)$  is (strict) anti-causal if and only if  $h(n)$  is zero for  $n > 0$  ( $n \geq 0$ ).*

For a general mixed causal, anti-causal LTI system with an impulse response  $h(n)$  that has a z-transform given as:

$$H(z) = \sum_{n=-\infty}^{\infty} h(n)z^{-n}$$

can be split into its causal part, denoted as  $[H(z)]_+$ , and strict anti-causal part, denoted as  $[H(z)]_-$ , as follows,

$$H(z) = \sum_{n=1}^{\infty} h(-n)z^n + \sum_{n=0}^{\infty} h(n)z^{-n} = [H(z)]_- + [H(z)]_+$$

Such a splitting may help determining the ROC of an impulse response as illustrated in the following example.

---

**Example 2.1 (ROC of mixed causal, anti-causal LTI system)**

Consider the following causal and strict anti-causal systems given by the z-transform of their impulse responses

$$H_c(z) = \frac{1}{1 - 0.9z^{-1}} \quad (\text{causal}) \quad H_a(z) = \frac{z}{1 - 0.9z} \quad (\text{strict anti-causal})$$

Then their series expansions are given resp. as,

$$H_c(z) = 1 + 0.9z^{-1} + 0.81z^{-2} + \dots + 0.9^n z^{-n} + \dots \quad H_a(z) = z + 0.9z^2 + 0.81z^3 + \dots + 0.9^n z^{n+1} + \dots$$

The ROCs of both z-transforms are:

$$\text{ROC}(H_c(z)) : \{z \in \mathbb{C} : |z| > 0.9\} \quad \text{ROC}(H_a(z)) : \{z \in \mathbb{C} : |z| < \frac{1}{0.9}\}$$

Therefore for the mixed causal, anti-causal system with z-transform of its impulse response given as,

$$H(z) = H_a(z) + H_c(z)$$

its ROC is the intersections of  $\text{ROC}(H_c(z))$  and  $\text{ROC}(H_a(z))$  and given as:

$$\text{ROC}(H(z)) : \{z \in \mathbb{C} : 0.9 < |z| < \frac{1}{0.9}\}$$


---

The splitting of a mixed causal, anti-causal into its causal part and its anti-causal part is done with a technique known as *partial fraction expansion*. The use of this technique is illustrated in the following example.

---

**Example 2.2 (Partial Fraction Expansion)**

Let the LTI system with transfer function  $H(z)$ :

$$\begin{aligned} H(z) &= \frac{1 - 1.8z^{-1} + z^{-2}}{-0.9 + 1.81z^{-1} - 0.9z^{-2}} \\ &= \frac{z - 1.8 + z^{-1}}{-0.9z + 1.81 - 0.9z^{-1}} \end{aligned}$$

This transfer function has poles 0.9 (stable pole) and  $\frac{1}{0.9}$  (unstable pole). We seek to split  $H(z)$  into a causal part and a strict anti-causal part with resp. transfer functions given as,

$$H_c(z) = \frac{A + Bz^{-1}}{1 - 0.9z^{-1}} \quad H_a(z) = \frac{Dz}{1 - 0.9z}$$

Then the coefficients  $A, B$  and  $D$  can be found from the following equation,

$$H(z) = \frac{A + Bz^{-1}}{1 - 0.9z^{-1}} + \frac{Dz}{1 - 0.9z} = \frac{Bz^{-1} + (A - 0.9B - 0.9D) + (-0.9A + D)z}{-0.9z^{-1} + 1.81 - 0.9z}$$

From this equation it follows,

$$A - 0.9B - 0.9D = -1.8 \quad -0.9A + D = 1 \quad B = 1$$

This yields  $A = 0, B = 1$  and  $D = 1$  or,

$$H(z) = \frac{z^{-1}}{1 - 0.9z^{-1}} + \frac{z}{1 - 0.9z}$$

Here we split the transfer function  $H(z)$  into stable causal and stable anti-causal parts with ROC  $0.9 < |z| < \frac{1}{0.9}$ , so the ROC includes the unit circle.

---

Our main interest will be in rational transfer functions that have the following form,

$$H(z) = \frac{b(0) + b(1)z^{-1} + \dots + b(q)z^{-q}}{1 + a(1)z^{-1} + \dots + a(p)z^{-p}} \quad (2.20)$$

We assume in general the numerator  $b(0) + b(1)z^{-1} + \dots + b(q)z^{-q}$  and the denominator  $1 + a(1)z^{-1} + \dots + a(p)z^{-p}$  to have no factors in common. This assumption is indicated by the numerator and denominator to be *coprime*. For a rational transfer function the roots of the numerator are called the *zeros* of the system and the roots of the denominator are called the *poles* of the system. Let  $p_i \in \mathbb{C}$  be the poles of the LTI system with transfer function  $H(z)$  as given in (2.20) such that its numerator and denominator are coprime, then this LTI is causal and BIBO if and only if,

$$|p_i| < 1 \quad (2.21)$$



With the  $z$ -transform as defined in (2.12), the corresponding time-domain difference equation representation of an LTI system can be derived from its transfer function when it is given as (2.20). Then by using the 'Time-shift' property of the  $z$ -transform in Table 2.3, the input-output representation given in (2.19) corresponds to the difference equation:

$$y(n) + a(1)y(n-1) + \cdots + a(p)y(n-p) = b(0)x(n) + b(1)x(n-1) + \cdots + b(q)x(n-q) \quad (2.22)$$

Here no initial conditions are considered. This is consequence by the use of the (double-sided) definition of the  $z$ -transform in (2.12). By the one-to-one correspondence of difference equations such as in (2.22) and the input-output representation given in (2.19) for  $H(z)$  given as in (2.20), the relationship between the time signals  $x(n)$  and  $y(n)$  will be graphically represented as in Figure 2.2.



**Figure 2.2:** Schematic Representation of a LTI Discrete-Time system represented by the transfer function  $H(z)$  (2.20) transforming its input  $x(n)$  to its output  $y(n)$  as governed by the difference equation (2.22).

**Definition 2.8** (Inverse of an LTI system). *The transfer function of the inverse of an LTI system given by its transfer function  $H(z)$  is:*

$$H^{-1}(z) = \left( H(z) \right)^{-1}$$

---

**Example 2.3 (Inverse of an LTI system)**

Consider the LTI system with transfer function:

$$H(z) = \frac{1}{1 - 0.9z^{-1}}$$

Then the transfer function of the inverse system is,

$$H^{-1}(z) = 1 - 0.9z^{-1} \quad (2.23)$$


---

**Definition 2.9** (Minimum-phase LTI system). *An LTI system with transfer function  $H(z)$  is minimum-phase if both the LTI system itself and its inverse are stable and causal.*

For LTI systems with rational transfer function, the minimum-phase property requires its poles and zeros to be within the unit disk, i.e their magnitude should be smaller than one.

The case when the degree of the numerator is smaller than that of the denominator, the checking of the minimum phase property requires the use of the original definition of stability. This is illustrated in the next example.

---

**Example 2.4 (Minimum-phase LTI system)**

Consider the LTI system with transfer function:

$$H(z) = \frac{1}{1 - 0.9z^{-1}}$$

This LTI system is stable since its pole is within the unit disk in the complex plane, i.e.  $0.9 < 1$ . The inverse system has transfer function,

$$H^{-1}(z) = 1 - 0.9z^{-1}$$

Following Lemma 2.6 the system with transfer function  $H^{-1}(z)$  is stable as well and the LTI system is minimum phase.

---

---

## 2.4 Optimizing cost functions w.r.t. a complex parameter

In solving inverse problems one often has to solve a parametric optimization problem as illustrated in the following example.

---

**Example 2.5 (A simple Inverse Problem)**

When the relationship between two (complex) signals  $x(n), y(n)$  is assumed to be governed by a first order difference equation,

$$y(n) - ay(n-1) = bx(n) \quad (2.24)$$

for  $a, b$  two unknown complex constants. Then a simple inverse problem is to deduce the values of the unknowns from a (finite) time record of the signals  $x(n), y(n)$ . Let us assume that the latter are known for  $n = 0 : N-1$ , an optimization problem to find the unknown constants can be stated as follows,

$$\min_{a, b \in \mathbb{C}} J(a, b) \quad \text{with} \quad J(a, b) = \sum_{k=1}^{N-1} \left| y(k) - [y(k-1) \quad x(k)] \begin{bmatrix} a \\ b \end{bmatrix} \right|^2 \quad (2.25)$$

---

The difficulty in optimizing cost functions like in (2.25) is that functional is not differentiable due to the presence of the modulus of the error  $e(k) = y(k) - [y(k-1) \quad x(k)] \begin{bmatrix} a \\ b \end{bmatrix}$ . The function  $|e(k)|^2$  depends on the conjugates  $a^*$  (and  $b^*$ ) and therefore is not differentiable with respect to  $a$  (or  $b$ ) [4]. An elegant way to resolve this problem is to treat the parameters  $a, b$  and their conjugate  $a^*, b^*$  as independent variables [5]. The use of this approach is again illustrated for the example 2.5

---

**Example 2.6 (A simple Inverse Problem (continuation of Example 2.5))**

We write the real-valued cost function  $J(a, b)$  in (2.25) compactly in terms of the error term  $e(k)$  as

$$e(k; a, b) = y(k) - [y(k-1) \ x(k)] \begin{bmatrix} a \\ b \end{bmatrix}$$

compactly as,

$$J(a, b) = \sum_{k=1}^{N-1} e(k; a, b) e^*(k, a, b)$$

Then following the approach outlined in [5], the necessary (and for this quadratic cost function also sufficient) condition for optimality is either given by:

$$\sum_{k=1}^{N-1} e(k; a, b) \begin{bmatrix} \frac{\partial e^*(k; a, b)}{\partial a} \\ \frac{\partial e^*(k; a, b)}{\partial b} \end{bmatrix} = 0 \quad (2.26)$$

or

$$\sum_{k=1}^{N-1} \begin{bmatrix} \frac{\partial e(k; a, b)}{\partial a} \\ \frac{\partial e(k; a, b)}{\partial b} \end{bmatrix} e^*(k; a, b) = 0 \quad (2.27)$$

This yields using (2.26) the following equations for the parameters  $a$  and  $b$ .

$$\sum_{k=1}^{N-1} \begin{bmatrix} y^*(k-1) \\ x^*(k) \end{bmatrix} \left( y(k) - [y(k-1) \ x(k)] \begin{bmatrix} a \\ b \end{bmatrix} \right) = 0$$

or written as an explicit set of equations known as the *Normal Equations* for problem (2.25),

$$\sum_{k=1}^{N-1} \begin{bmatrix} y^*(k-1) \\ x^*(k) \end{bmatrix} y(k) = \sum_{k=1}^{N-1} \begin{bmatrix} y^*(k-1) \\ x^*(k) \end{bmatrix} [y(k-1) \ x(k)] \begin{bmatrix} a \\ b \end{bmatrix} \quad (2.28)$$

---

---

## References

- [1] S. Mitra, *Digital Signal Processing*. New York: McGraw-Hill, 1998.
- [2] A. Oppenheim and R. Schaffer, *Discrete-Time Signal Processing*. New Jersey: Prentice Hall, 2nd ed., 1998.
- [3] C.-T. Chen, *Signals and Systems*. New York Oxford: Oxford University Press, third ed., 2004.
- [4] M. H. Hayes, *Statistical Digital Signal Processing and Modeling*. New York: John Wiley and Sons, 1996.

- [5] D. Brandwood, "A complex gradient operator and its application in adaptive array theory," *IEE Proc. Parts F and H*, vol. 130, pp. 11–16, 1983.
- [6] E.W.Kamen, *Introduction to Signals and Systems*, Macmillan Publishing Company, 1990.
- [7] W.R. LePage, *Complex Variables and the Laplace Transform for Engineers*, Dover Publications, 1980.

## Exercises

**Exercise 2.1** Let the signal  $x(n)$  be represented as in (2.2) then use the linearity property of the DTFT in Table 2.2 to show that,

$$\mathcal{F}(x(n)) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n}$$

**Exercise 2.2** Prove that the following relationship holds,

$$\sum_{n=-\infty}^{\infty} |x(n)|^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |X(e^{j\omega})|^2 d\omega \quad (2.29)$$

This relationship is known as *Parseval's Theorem*.

**Exercise 2.3** When a signal  $x(n)$  is conjugate symmetric (as we encounter in the analysis of signal spectra in Chapter 5), show that its z-transform has the following property,

$$X(z) = X^*\left(\frac{1}{z^*}\right) \quad (2.30)$$

**Exercise 2.4** Determine the z-transform and its ROC of the unit step  $u(n)$  defined in Table 2.1.

**Exercise 2.5** Consider the convolution in (2.17), then

- (a) Show that for general mixed causal, anti-causal systems (2.17) can also be expressed as,

$$y(n) = \sum_{k=-\infty}^{\infty} h(k)x(n-k) \quad (2.31)$$

- (b) Show that when the LTI system with impulse response  $h(n)$  is causal, (2.17) can be written as,

$$y(n) = \sum_{k=0}^{\infty} h(k)x(n-k) \quad (2.32)$$

**Exercise 2.6** Let the following transfer function  $H(z)$  be given, with the unit circle in its ROC:

$$H(z) = \frac{(1 - 0.8z^{-1})(1 - 2z^{-1})}{(1 - 0.7z^{-1})(1 - 0.9z^{-1})}$$

then,

- (a) Determine whether or not  $H(z)$  is causal, and justify your answer.
- (b) Determine whether or not  $H(z)$  is stable, and justify your answer.
- (c) Determine whether or not  $H(z)$  is minimum-phase, and justify your answer.
- (d) Split the inverse  $H^{-1}(z)$  into the sum of a causal and an anti-causal part.



# Chapter 3

## Random Variables

**After studying this chapter you can:**

- characterize a single (or a pair of) random variable(s) via its (their joint) Probability Distribution function or Probability Density (pdf) function,
- use the statistical notions of independency, uncorrelatedness and orthogonality of two random variables and you know what the pdf of a single or a pair of Gaussian random variables is.

---

## 3.1 Introduction

In this chapter we start with a review of the description of random variables. This is relevant since stochastic processes are (ordered) sequences of random variables. The general description of such sequences would require the definition of the Probability Density function of all individual elements in that sequence and of all their mutual combinations. This in general is extremely complicated and for the purpose of this course we make two restrictions. First we will describe stochastic processes only in terms of their so-called first order and second order statistical moments. The first order statistical moment being the mean and the second order the correlation or the covariance function. In general, unless otherwise stipulated, zero-mean stochastic processes will be considered. The second restriction is that we will assume the stochastic processes to be stationary. After a definition of this notion the important concept of Wide Sense Stationary (WSS) is introduced. In this course we will mainly be focusing on the analysis and synthesis of zero-mean WSS stochastic processes. The notion of ergodicity is introduced and its use is illustrated for the estimation of the mean value of a WSS stochastic process.

In Section 5.2.8 spectra are defined as Fourier transforms of the second order statistical moments of zero-mean WSS stochastic processes. An introduction is made on how to derive these spectra from a single and finite length realization of the stochastic process.

---

## 3.2 Random Variables

A random variable is the outcome of an experiment, e.g. throwing dice. An example of a physical random variable is given first.

---

### Example 3.1 (Example 1.2 (Ct'd))

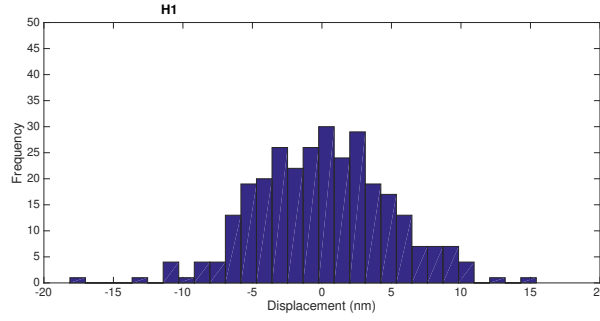
Please consider the displacement of a particle “suspended” in a medium as displayed in Figure 1.2 when repeating the experiment multiple times. Then the position of a particle (e.g. in the  $x$ -direction at **one** time instance (e.g.  $0.05\mu s$ ) can be represented as a *continuous random variable*.

We can represent the “distribution” of the position at the time instance  $0.05\mu s$  with a histogram [1]. Such an histogram is displayed in Figure 3.1 for 300 repeated (simulation) experiments. For the data used in the simulation we refer to Example 1.2.

---

In Example 3.1 the random variable is the  $x$ -position of a particle in a water solution at time instance  $0.05\mu s$ . This random variable is the outcome of a random experiment. In the case of Example 3.1 the experiment is the observation (with a microscope) of a particle moving in a water solution. What makes the experiment random is that each time we repeat the experiment with the particle in





**Figure 3.1:** A histogram of the  $x$ -displacement (nm) of a particle governed by the Langevin Equation (1.9) with data provided in Example 1.2 at time instance  $0.05\mu\text{s}$  and repeating the experiment 300 times starting from identical initial conditions.

the same initial position, a different evolution of the particle position over time results.

The outcome of random (“non-repeatable”) experiments need, contrary to example 3.1 not to be continuous. It can be discrete and it can even be not a number. This is discussed in the next two subsections.

### 3.2.1 Discrete Random Variables

A discrete random variable with outcomes not being numbers is the flipping of a “fair” coin [4]. A “fair” coin is equally likely to give Heads or Tails. In this case the *Space of all possible Outcomes*, which is also called the *sample space*, consists of the pair {Heads, Tails}. The sample space is denoted by the symbol  $\Omega$  and in the case of flipping a coin equals:

$$\Omega = \{\text{Heads, Tails}\}$$

In order to solely work with numbers we can assign to each possible outcome a real number. For example we assign to the element “Heads” of  $\Omega$  the real number +1 and to the element “Tails” the number  $-1$ . An element or a subset of elements of the set  $\Omega$  is called *an event*. An event may even be a complex number, such as for the random variable,

$$z = a + bj$$

with  $a$  the outcome of flipping a coin and  $b$  the outcome of throwing a die.

A random variable is fully characterized if we assign to each entry of  $\Omega$  a *probability*. If we denote this assignment by the operator  $\text{Pr}(\cdot)$ , then for a fair coin we have:

$$\text{Pr}(\text{Heads}) = 0.5 \quad \text{Pr}(\text{Tails}) = 0.5$$

The Histogram of the outcome of flipping a “fair” coin would provide for the two elements of  $\Omega$  a frequency value that indeed approaches 0.5 for both random variables “Heads” and “Tails” for a large number of “trials” with a fair coin. Let  $\#H$  denote the number of times Heads was observed in an experiment flipping a fair coin  $N$  times and let  $\#T$  denote the number of Tails, then the frequency for the event “Heads” equals  $\frac{\#H}{N}$ , while that for the event “Tails” equals  $\frac{\#T}{N}$ .

**Definition 3.1** (Discrete Random Variable). *A discrete Random Variable with sample space  $\Omega = \{\omega_i\}$  for  $i = 1 : n$  is fully characterized if we assign a probability  $p_i$  to each element of  $\Omega$  following the axioms given in Table 3.1, denoted as,*

$$\Pr(\omega_i) = p_i \in [0, 1] \quad (3.1)$$

|          |   |
|----------|---|
| Axiom 1: | For any event $A \in \Omega$ : $\Pr(A) \geq 0$  |
| Axiom 2: | For the certain event $\Omega$ : $\Pr(\Omega) = 1$  |
| Axiom 3: | For any two events $A_1$ and $A_2$ in $\Omega$ :<br>$\Pr(A_1 \cup A_2) = \Pr(A_1) + \Pr(A_2) - \Pr(A_1 \cap A_2)$ |

**Table 3.1:** The axioms of assigning a probability to the outcomes of a (discrete or continuous) random variable. Here  $\Omega$  denotes the sample space of a random variable (all possible outcomes) and  $A_1$  and  $A_2$  denote different events (subsets of outcomes).

**Definition 3.2** (Bernoulli random variable). *When we have a random variable with two possible outcomes, e.g. numerically indicated by  $\omega_1 = 1$  and  $\omega_2 = -1$ . Then these outcomes represent a Bernoulli random variable if we assign the following probabilities to each possible outcome (event);*

$$\Pr(\omega_1) = p \in [0, 1] \quad \Pr(\omega_2) = 1 - p$$

Which we could also denote as,

$$\Pr(\omega = 1) = p \in [0, 1] \quad \Pr(\omega = -1) = 1 - p$$

An example of a Bernoulli random variable is the flipping of an “unfair” coin where one side is more likely to occur than the other side of the coin.

### 3.2.2 Continuous Random Variables

The outcome of the  $x$ -position e.g. at the time instance  $0.05\mu s$  as defined in the Example 3.1, is a real variable in  $\mathbb{R}^+$ . The sample space in this case is

$$\Omega = \{x : x \in \mathbb{R}^+\}$$

Such a random variable is called *continuous* for the outcome is a continuum of values. For a continuous random variable, we assign probabilities to subsets of  $\Omega$  rather than to individual elements. Such assignment should be done taken the axioms given in Table 3.1 into account.

The subsets are subintervals in  $\mathbb{R}^+$ . Considering Example 3.1 the position of a particle at  $0.05\mu s$ , a subset could be the interval  $(5, 10]$  nm.

### 3.2.3 Characterizing Random Variables

#### The Probability Distribution and Density Function

The random variable is fully characterized by assigning probabilities to its outcomes. Such assignment can be done via the *Probability Distribution Function (PDF)*. If we denote the random variable by  $x$  then its PDF is a function that assigns a probability to the event  $\{x \leq \alpha\}$ , denoted as:

$$F_x(\alpha) = \Pr\{x \leq \alpha\} \quad (3.2)$$

Though a probability is assigned to an interval of the sample space, it will be seen that the PDF can be used for both discrete and continuous random variables. This is illustrated in Example 3.3.

In practice use is made of the *probability density function (pdf)* to characterize a random variable. Let for a random variable  $x$  the pdf be denoted by the function  $f_x(\alpha)$  then this function is the derivative of the PDF of  $x$ :

$$f_x(\alpha) = \frac{dF_x(\alpha)}{d\alpha} \quad (3.3)$$

---

#### Example 3.2 (Uniform continuous Random Variable)

Consider the protective coating on a conductor so to protect the conductor to operate in a corrosive environment. When the thickness  $t$  of the protective layer is equally likely to occur in the interval  $[0, 1]$  ( $= \Omega$ ), then the probability assignment to the random variable  $t$  can be done as follows:

$$\text{Let } \mathcal{T} = (\alpha_1, \alpha_2] \subset \Omega \text{ then } \Pr\{t \in \mathcal{T}\} = \Pr\{\alpha_1 < t \leq \alpha_2\} = \alpha_2 - \alpha_1$$

That is the probability of an event, i.e. any subinterval in  $\Omega$ , depends on the size of that subinterval.

For the random variable  $t$  the PDF is defined as follows:

$$F_t(\alpha) = \begin{cases} 0 & : \alpha < 0 \\ \alpha & : 0 \leq \alpha < 1 \\ 1 & : \alpha \geq 1 \end{cases} \quad (3.4)$$

The PDF and its derived pdf function of the uniform random variable  $t$  is displayed in Figure 3.2.

---

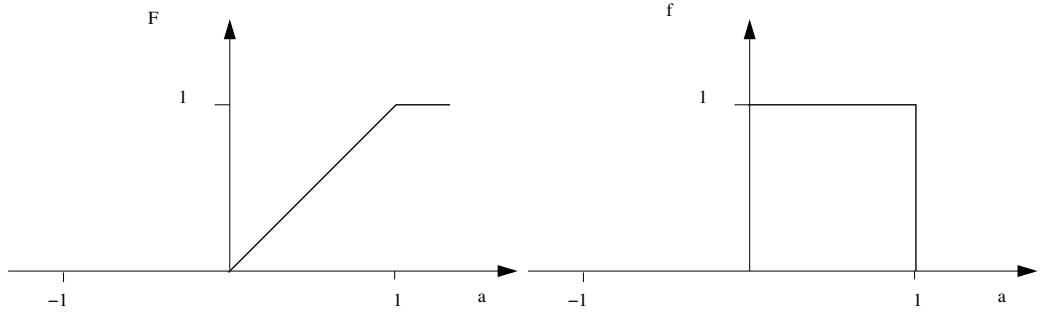
#### Example 3.3 (Uniform discrete Random Variable)

In this example we consider a “fair” die. Then the random variable is the number of eyes  $e$  we get after throwing the die. In this case the sample space equals:

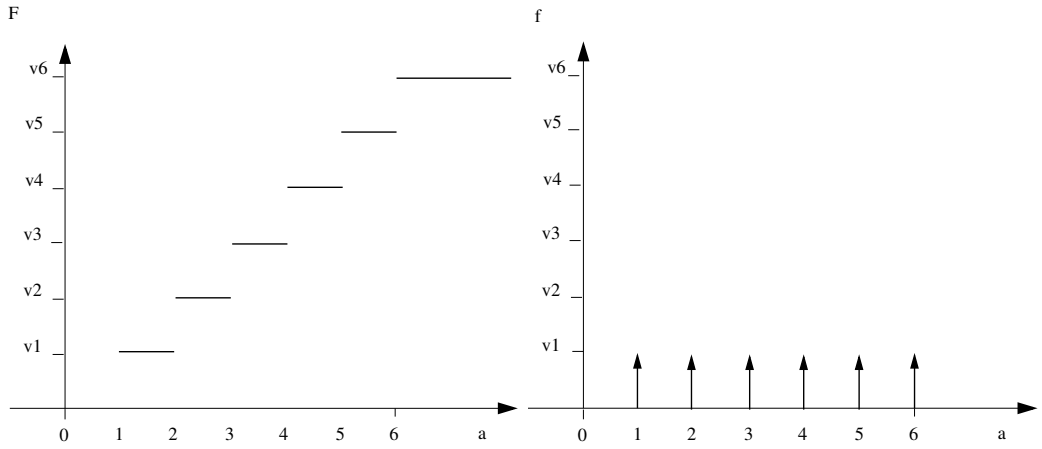
$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

When each element in  $\Omega$  is equally likely to occur, the probability of each entry is  $\frac{1}{6}$  and we obtain the PDF and pdf as displayed in Figure 3.3.

---



**Figure 3.2:** The Probability Distribution Function (PDF) and its derivative (probability density function - pdf) of a uniform continuous random variable.



**Figure 3.3:** The Probability Distribution Function (PDF) and its derivative (probability density function - pdf) of a uniform discrete random variable.

### Ensemble Averages

In practice it is generally sufficient to characterize an average behavior of a random variable. For example consider Example 3.1, rather than the full characterization of the distribution of the  $x$ -position at  $0.05\mu$  s via its PDF or pdf, one could be interested in the mean value of the square of the displacement. The latter question was a topic of interest to the thesis of A. Einstein [?].

For a definition of the mean or the *expected value* or the *first moment* of a random variable use is made of the so-called *Expectation Operator* denoted as  $E[.]$ . For a discrete random variable  $x$  with  $\Omega$  defined as:

$$\Omega = \{x_1, x_2, \dots, x_n\} \quad (3.5)$$

the mean value is given as:

$$E[x] = \sum_{k=1}^n x_k \Pr\{x = x_k\} \quad (3.6)$$

For a continuous random variable  $x$  with  $\Omega$  the entire real axis and with pdf

denoted by  $f_x(\alpha)$ , the mean value is given as:

$$E[x] = \int_{-\infty}^{\infty} \alpha f_x(\alpha) d\alpha \quad (3.7)$$

The resemblance between the definition of the mean for a continuous random variable in (3.7) with that for a discrete random variable in (3.6) indicates that the ‘measure’  $f_x(\alpha)d\alpha$  could be ‘interpreted’ as probability assigned to the event  $\alpha$ .

**Remark 3.3.** When we denote the pdf of a discrete random variable  $x$  with sample space  $\Omega$  in (3.5) as:

$$f_x(\alpha) = \sum_{k=1}^n \Pr\{x = x_k\} \delta(\alpha - x_k)$$

it is left as an exercise to prove that (3.7) equals (3.6). This definition of the pdf would allow us to just use the formulas of the ensemble averages defined next for continuous random variables for discrete random variables as well.

---

**Example 3.4 (Determining the mean of a Random Variable given its full probabilistic characterization)**

Consider the random variable  $e$  in Example 3.3, then according to (3.6) we have:

$$E[e] = \sum_{k=1}^6 k \Pr\{e = k\} = 3.5$$

It should be remarked that in this case the mean value does not belong to the sample space  $\Omega$ . That is the mean value will never occur.

For the continuous random variable  $t$  in Example 3.2, the mean value is:

$$E[t] = \int_0^1 \alpha f_t(\alpha) d\alpha = \int_0^1 \alpha d\alpha = 0.5$$

In this case since  $t$  is continuous, this value may occur as an outcome of an experiment.

---

A generalization of the mean of a random variable could be the mean of a nonlinear transformation of a random variable. One widely used generalization is the second moment of a random variable or the *mean square value*. For a discrete random variable with sample space in (3.5) it is defined as,

$$E[x^2] = \sum_{k=1}^n x_k^2 \Pr\{x = x_k\} \quad (3.8)$$

For a continuous random variable with pdf  $f_x(\alpha)$  and sample space the entire real axis, it is defined as:

$$E[x^2] = \int_{-\infty}^{\infty} \alpha^2 f_x(\alpha) d\alpha \quad (3.9)$$

A modification of the second moment is to consider the average of the squared value of the difference between the random variable and its mean. This is called the *variance* of a random variable. It usually denoted by the symbol  $\sigma_x^2$  and is defined as:

$$\sigma_x^2 = E[(x - E[x])^2] \quad (3.10)$$

When we have an estimate of the random variable  $x$  denoted as  $\hat{x}$ , the *Mean Square Error* is defined as,

$$MSE = E[(x - \hat{x})^2] \quad (3.11)$$

When  $z$  is a complex random variable its variance is real and is defined as:

$$\sigma_z^2 = E[(z - E[z])(z^* - E[z]^*)]$$

In a similar way the mean of a nonlinear transformation of a random variable  $x$  can be defined. Let  $x$  and  $g(x)$  be random variables such that,

$$y = g(x)$$

then the expected value of  $y$  is,

$$E[y] = E[g(x)] = \int_{-\infty}^{\infty} g(\alpha) f_x(\alpha) d\alpha \quad (3.12)$$

### 3.2.4 Two Random Variables

Two random variables are the outcomes of two random processes. For example throwing 2 dices at ones or measuring the  $x$  and  $y$  position of a particle in a water solution. Then we could generalize the concepts used to describe one random variables to two or more. The latter is necessary in the analysis of stochastic processes.

---

#### Example 3.5 (Two “fair” coins)

When flipping two fair coins, and when indicating “Heads” by +1 and “Tails” by -1, then we can denote the sample space  $\Omega$  as:

$$\Omega = \{(-1, -1), (1, -1), (-1, 1), (1, 1)\}$$

and the probability that can be assigned to each entry of  $\Omega$  is  $\frac{1}{4} = \frac{1}{2} \frac{1}{2}$ .

---

### The Joint Distribution and Density Function

For two random variables  $x(1)$  and  $x(2)$  we can define the *joint distribution function* as:

$$F_{x(1),x(2)}(\alpha_1, \alpha_2) = Pr(x(1) \leq \alpha_1 \text{ and } x(2) \leq \alpha_2) \quad (3.13)$$

and the *joint density function* as:

$$f_{x(1),x(2)}(\alpha_1, \alpha_2) = \frac{\partial^2}{\partial \alpha_1 \partial \alpha_2} F_{x(1),x(2)}(\alpha_1, \alpha_2) \quad (3.14)$$

## Joint Moments

As a generalization of the ensemble averages defined for one random variable we can also define ensemble averages for two random variables. The overview of the definition of these moments is given Table 3.2. In this table it is assumed that the two continuous, random variables  $x$  and  $y$  have joint density function:

$$f_{x,y}(\alpha, \beta)$$

and that their sample space  $\Omega = \mathbb{R} \times \mathbb{R}$ .

| Quantity                 | Symbol       | Definition  |
|--------------------------|--------------|---|
| Means:                   | $E[x]$       | $E[x] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \alpha f_{x,y}(\alpha, \beta) d\alpha d\beta$   |
|                          | $E[y]$       | $E[y] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \beta f_{x,y}(\alpha, \beta) d\alpha d\beta$  |
| Variances:               | $\sigma_x^2$ | $\sigma_x^2 = E[(x - E[x])^2]$<br>$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\alpha - E[x])^2 f_{x,y}(\alpha, \beta) d\alpha d\beta$ |
|                          | $\sigma_y^2$ | $\sigma_y^2 = E[(y - E[y])^2]$<br>$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\beta - E[y])^2 f_{x,y}(\alpha, \beta) d\alpha d\beta$  |
| Correlation:             | $r_{xy}$     | $E[xy] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \alpha \beta f_{x,y}(\alpha, \beta) d\alpha d\beta$                                  |
| Covariance:              | $c_{xy}$     | $E[(x - E[x])(y - E[y])]$   |
| Correlation Coefficient: | $\rho_{xy}$  | $\frac{E[(x - E[x])(y - E[y])]}{\sigma_x \sigma_y}$   |

**Table 3.2:** Different Moments for two continuous random variables with joint density function  $f_{x,y}(\alpha, \beta)$  and sample space  $\Omega = \mathbb{R} \times \mathbb{R}$ .

The joint moments defined in Table 3.2 can also be defined for pairs of complex random variables [4]. For example let  $x$  and  $y$  be a pair of continuous, complex random variables then their correlation equals:

$$r_{xy} = E[xy^*] \quad (3.15)$$

The definition of the expectation operator for two random variables as done in Table 3.2 enables to show a key property of this operator, namely that it is a *linear operator*. This is demonstrated in the following Lemma.

**Lemma 3.4.** When  $x$  and  $y$  are two random variables with sample space the product space  $\mathbb{R} \times \mathbb{R}$  and with joint pdf denote by  $f_{x,y}(\alpha, \beta)$ , and let  $a, b \in \mathbb{C}$  be two constants, then

$$E[ax + by] = aE[x] + bE[y] \quad (3.16)$$

*Proof.* Using the definition of the expectation of two random variables, and the facts that  $\int_{-\infty}^{\infty} f_{x,y}(\alpha, \beta) d\alpha = f_y(\beta)$  and  $\int_{-\infty}^{\infty} f_{x,y}(\alpha, \beta) d\beta = f_x(\alpha)$  we have that:

$$\begin{aligned}
 E[ax + by] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (a\alpha + b\beta) f_{x,y}(\alpha, \beta) d\alpha d\beta \\
 &= a \int_{-\infty}^{\infty} \alpha \left( \int_{-\infty}^{\infty} f_{x,y}(\alpha, \beta) d\beta \right) d\alpha + b \int_{-\infty}^{\infty} \beta \left( \int_{-\infty}^{\infty} f_{x,y}(\alpha, \beta) d\alpha \right) d\beta \\
 &= a \int_{-\infty}^{\infty} \alpha f_x(\alpha) d\alpha + b \int_{-\infty}^{\infty} \beta f_y(\beta) d\beta \\
 &= aE[x] + bE[y]
 \end{aligned}$$

□

We end this section by postulating the important *Cauchy Schwartz* inequality of two random variables.

**Lemma 3.5** (Cauchy-Schwartz inequality). [4] Let  $x$  and  $y$  be two random variables with joint density function:

$$f_{x,y}(\alpha, \beta)$$

and sample space  $\Omega = \mathbb{R} \times \mathbb{R}$ , let the covariance, variance of  $x$  and variance of  $y$  be resp. denoted as in Table 3.2, then,

$$|c_{xy}|^2 \leq \sigma_x^2 \sigma_y^2$$

*Proof.* A proof for the special case that  $x$  and  $y$  have mean zero is called for in Exercise 3.7. □

### Independence, Uncorrelatedness and Orthogonality

Two random variable may be related or unrelated. For example when you throw 2 coins, which are fair, the outcome of one coin has nothing to do with the outcome of the other coin. This may not be the case between the  $x$  and  $y$  position of a particle moving in a water solution.

The strongest notion of two random variables not being related to one another is *independence*. This notion is defined next.

**Definition 3.6** (Independence between 2 random variables). Two random variables with joint density function  $f_{x,y}(\alpha, \beta)$  and individual density functions resp. given as  $f_x(\alpha)$ ,  $f_y(\beta)$  are statistically independent if,

$$f_{x,y}(\alpha, \beta) = f_x(\alpha) f_y(\beta)$$

A weaker form of expressing that two random variables are not related is the notion of uncorrelated random variables. In that case we have,

$$E[xy^*] = E[x]E[y^*] \quad (3.17)$$

Finally two random variables are *orthogonal* if we have

$$E[xy^*] = 0. \quad (3.18)$$

From these definitions it is clear that two random variables that are uncorrelated, they are also orthogonal provided one of the random variables has mean zero.

### 3.2.5 Gaussian Random Variables

#### For one Random Variable

When  $x$  is a continuous random variable with a Gaussian distribution, for short a Gaussian Random Variable, and moments as defined in Table 3.2, its pdf is given by:

$$f_x(\alpha) = \frac{1}{\sigma_x \sqrt{2\pi}} e^{-\frac{(\alpha - E[x])^2}{2\sigma_x^2}} \quad (3.19)$$



This distribution function allows us to determine the (symmetric) interval around the mean value of  $x$  such that the outcome  $x$  is within this interval with a certain probability. For example with probability  $p = 0.682689492137 \dots$  (about 68%) of the values of  $x$  lie in an interval  $[E[x] - \sigma_x, E[x] + \sigma_x]$ .

The above value of  $p$  is retrieved from the  $\chi^2$  distribution with one degree [?].

### For two Random Variables

When  $x$  and  $y$  are continuous random variables with a Gaussian distribution, and moments as defined in Table 3.2, we can define their covariance matrix  $C$  as:

$$C = E \left[ \begin{pmatrix} x - E[x] \\ y - E[y] \end{pmatrix} \begin{pmatrix} x - E[x] & y - E[y] \end{pmatrix} \right] = \begin{bmatrix} \sigma_x^2 & c_{xy} \\ c_{xy} & \sigma_y^2 \end{bmatrix} \quad (3.20)$$

Then the joint density function is given by:

$$f_{x,y}(\alpha, \beta) = \frac{1}{2\pi\sqrt{|\det C|}} e^{-\frac{1}{2} \begin{bmatrix} \alpha - E[x] & \beta - E[y] \end{bmatrix} C^{-1} \begin{bmatrix} \alpha - E[x] \\ \beta - E[y] \end{bmatrix}} \quad (3.21)$$

For the case the covariance matrix  $C$  is the identity matrix, again the  $\chi^2$  distribution with degree two provides a way to determine the radius of the circle with the mean vector  $\begin{bmatrix} E[x] \\ E[y] \end{bmatrix}$  as center, that contains the pairs  $(x, y)$  in the 2D Euclidean plane with a certain probability. When the covariance matrix  $C$  is different from the identity matrix the distance to the mean vector is weighted in some sense.

---

#### Example 3.6 (2D Gaussian Random Variable)

Let us consider two zero-mean, random variables  $x$  and  $y$  that have a Gaussian distribution. When their covariance matrix is taken as:

$$C = \begin{bmatrix} \frac{1}{4} & .3 \\ .3 & 1 \end{bmatrix}$$

the joint density function is displayed in Figure 3.4.

---



---

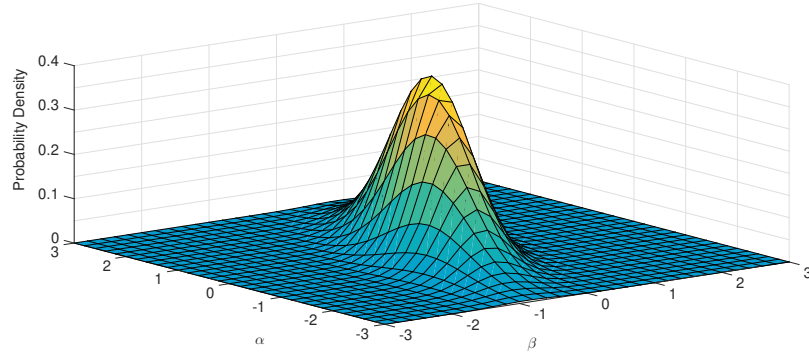
## Exercises

**Exercise 3.1** When  $x$  is a random variable with pdf given as:

$$f_x(\alpha) = \frac{1}{\sqrt{2\pi}} e^{-\frac{\alpha^2}{2}}$$

then determine the PDF of  $x$ .

**Exercise 3.2** Compute the variance of the random variables in Example 3.4.



**Figure 3.4:** The joint (probability) density function of 2 Gaussian distributed random variables with mean zero.

**Exercise 3.3** Determine for which value of the mean of a random variable  $x$  the second moment equals the variance of a random variable.

**Exercise 3.4** Compute the second moment and the variance of the random variable with pdf  $f_x(\alpha)$  given as:

$$f_x(\alpha) = \begin{cases} 1 & : 0 \leq \alpha \leq 1 \\ 0 & : \text{otherwise} \end{cases}$$

**Exercise 3.5** Prove the statement in Remark 3.3.

**Exercise 3.6** Let  $x$  be a random variable with sample space  $\Omega = \mathbb{R}$  and with pdf given as  $f_x(\alpha)$  then prove that,

$$\int_{-\infty}^{\infty} f_x(\alpha) d\alpha = 1$$

**Exercise 3.7** Prove the Cauchy-Schwartz inequality for the case  $x$  and  $y$  in addition to the conditions stated in Lemma 3.5 both have mean zero.

[ **Hint:** Start from the fact that for any real number  $\alpha$  we have that,

$$E[(\alpha x - y)^2] \geq 0 \quad ]$$

**Exercise 3.8** Prove that for the moments defined in Table 3.2 the correlation coefficient between 2 random variables  $x$  and  $y$  satisfies:

$$|\rho_{xy}| \leq 1$$

**Exercise 3.9** Prove that for the moments defined in Table 3.2 and considering the case of complex random variables, we have that,

$$c_{xy} = r_{xy} - E[x]E[y^*]$$

**Exercise 3.10** When  $x$  and  $y$  are two uncorrelated, zero-mean, complex random variables, prove that

$$E[|x + y|^2] = E[|x|^2] + E[|y|^2]$$

(in other words, the variance of the sum of two uncorrelated random variables is the sum of their variances).



# Chapter 4

## Estimation

**After studying this chapter you can:**

- explain the basic concepts for estimating unknown physical quantities on the basis of (random) measurement data
- assess the quality of estimators in terms of bias and variance
- explain basic estimation principles as least-squares, linear regression and maximum likelihood
- understand and derive the Cramer-Rao lower bound, and its relation to linear regression and maximum likelihood for estimation problems with certain properties.

---

## 4.1 Introduction

In estimation theory it is the objective to determine/estimate an unknown physical quantity or variable, on the basis of data that is available from measurements, which are generally subject to stochastic uncertainties. Because of the stochastic uncertainties, the measurement data is actually described by a set (or sequence) of random variables.

An estimator (or sometimes called *statistic*)<sup>1</sup> is any (deterministic) function of a (set of) random variable(s); the estimator itself does not contain any unknown parameters.

Let, e.g.,  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be a random sample from the probability density function  $f_{\mathbf{x}}(x)$ . Then we can write:

$$\mathbf{y} = g(\mathbf{x}_1, \dots, \mathbf{x}_N)$$

where  $g$  is any function (an estimator) that gives an estimate  $\mathbf{y}$ . In other words, an estimator is simply a function of random variables, and attempts to reproduce/estimate an unknown (physical) quantity on the basis of measured data.

$$\bar{\mathbf{x}} := \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

is an estimator of the mean value of  $\mathbf{x}$ , but also

$$\tilde{\mathbf{x}} := \frac{1}{2} [\max_i \mathbf{x}_i + \min_i \mathbf{x}_i]$$

can be an estimator of this quantity.

Estimators are not bound to have any direct relation with the quantity/variable that is estimated. Even  $\min_i \mathbf{x}_i$  can be an estimator of the mean value of  $\mathbf{x}$ , although one may guess that the properties of this estimator will be rather poor, unless  $\mathbf{x}_i$  is fixed for all values of  $i$ .

In this chapter several approaches will be discussed on how to design estimators for particular objects and with particular properties. The discussion is directed towards general estimators and their properties. Estimation of second order moments (correlation functions and spectral density function) of stochastic processes will be the subject of chapter 5.

### 4.1.1 Basic definitions

We will first discuss several notions that characterize the “quality” of estimators. As a notation we will consider an estimator  $\hat{\theta}_N$  of an underlying quantity  $\theta_0$ , based on  $N$  measurements. Since we will generally assume that both  $\hat{\theta}_N$  and  $\theta_0$  are real-valued vectors in  $\mathbb{R}^d$ , the definitions will be given for real-valued

---

<sup>1</sup>A statistic is any deterministic function of observable random variables, which is itself an observable random variable, which does not contain any unknown parameters (Mood et al., 1974). When a statistic is used to estimate a function  $\tau(\theta)$  of an unknown parameter  $\theta$ , it is called an estimator of  $\tau(\theta)$ .

parameters. However, they can easily be extended to include complex-valued parameters.

### Bias

The estimator  $\hat{\theta}_N$  is called **unbiased** if

$$E[\hat{\theta}_N] = \theta_0,$$

i.e., the estimator is delivering the “right” quantity “on average” and for a particular finite value of  $N$ . Additionally the estimator is called **asymptotically unbiased** if

$$\lim_{N \rightarrow \infty} E[\hat{\theta}_N] = \theta_0.$$

Estimators that are not unbiased are called biased, and the bias is given by

$$bias(\hat{\theta}_N) = E[\hat{\theta}_N] - \theta_0.$$

### Variance

The variance of  $\hat{\theta}_N$  is given by the mean deviation from its expected value, measured in a quadratic sense, i.e.,

$$var(\hat{\theta}_N) = E \left[ \left( \hat{\theta}_N - E[\hat{\theta}_N] \right) \left( \hat{\theta}_N - E[\hat{\theta}_N] \right)^T \right].$$

It determines the variation of outcomes of  $\hat{\theta}_N$  about its mean value. For vector-valued parameters the resulting matrix is referred to as the covariance matrix, also denoted as

$$cov(\hat{\theta}_N)$$

having dimensions  $d \times d$  when  $\hat{\theta}_N$  has dimension  $d$ . For scalar-valued parameters this covariance matrix reduces to the variance of the random variable  $\hat{\theta}_N$ .

### Mean squared error (MSE)

The mean squared error (MSE) is given by<sup>2</sup>

$$MSE(\hat{\theta}_N) := E[(\hat{\theta}_N - \theta_0)^2].$$

Denote  $E[\hat{\theta}_N] = m$ , then

$$\begin{aligned} MSE(\hat{\theta}_N) &= E[(\hat{\theta}_N - m + m - \theta_0)^2] \\ &= E[(\hat{\theta}_N - m)^2] + (m - \theta_0)^2 + 2E[(\hat{\theta}_N - m)(m - \theta_0)]. \end{aligned}$$

Because  $E[\hat{\theta}_N] = m$  the last term on the right hand side is 0, and therefore

$$MSE(\hat{\theta}_N) = var(\hat{\theta}_N) + [bias(\hat{\theta}_N)]^2. \quad (4.1)$$

---

<sup>2</sup>We first consider the situation of scalar real-valued parameters.

This renders the measure  $\text{MSE}(\hat{\theta}_N)$  into a notion that includes aspects of both bias and variance. It describes the mean deviation of the estimator from its exact value, measured in a quadratic sense.

For vector-valued parameters a natural extension exists, writing

$$\text{MSE}(\hat{\theta}_N) := \text{trace}[E[(\hat{\theta}_N - \theta_0)(\hat{\theta}_N - \theta_0)^T]] = E[(\hat{\theta}_N - \theta_0)^T(\hat{\theta}_N - \theta_0)]$$

then through a similar reasoning as above,

$$\begin{aligned} \text{MSE}(\hat{\theta}_N) &= \text{trace}[E\left[(\hat{\theta}_N - m)(\hat{\theta}_N - m)^T\right] + E\left[(m - \theta_0)(m - \theta_0)^T\right]] \\ &= \text{trace}[\text{cov}(\hat{\theta}_N) + \text{bias}(\hat{\theta}_N) \cdot \text{bias}(\hat{\theta}_N)^T] \end{aligned} \quad (4.2)$$

$$= \sum_i \text{var}(\hat{\theta}_N^{(i)}) + \sum_i \text{bias}(\hat{\theta}_N^{(i)})^2 \quad (4.3)$$

where  $(\cdot)^{(i)}$  refers to the  $i$ -th component of a vector. As a result the MSE for a vector-valued parameter is equal to the sum of the MSE's of each of its components.

### Consistency

An estimator  $\hat{\theta}_N$  is called (weakly) consistent if for every  $\delta > 0$ ,

$$\lim_{N \rightarrow \infty} \Pr[\|\hat{\theta}_N - \theta_0\| > \delta] = 0$$

also denoted as  $\text{plim}_{N \rightarrow \infty} \hat{\theta}_N = \theta_0$ .

To test if an estimator is consistent the estimator needs to converge in probability. It means that for every  $\theta$  with  $\|\theta_0 - \theta\| > \delta$  with  $\delta$  arbitrarily small, the probability density function  $f_{\hat{\theta}_N}(\theta)$  disappears for  $N \rightarrow \infty$ . particular when the probability density function converges to a dirac-pulse for  $N \rightarrow \infty$ .

Note that a (weakly) consistent estimator is not necessarily unbiased for finite values of  $N$ , as is also illustrated in figure 4.1. However an estimator that is asymptotically unbiased and its covariance is asymptotically going to zero is consistent. This situation is illustrated in figure 4.1.

### Efficiency

An unbiased estimator  $\hat{\theta}_N$  is called an efficient estimator of  $\theta_0$  if

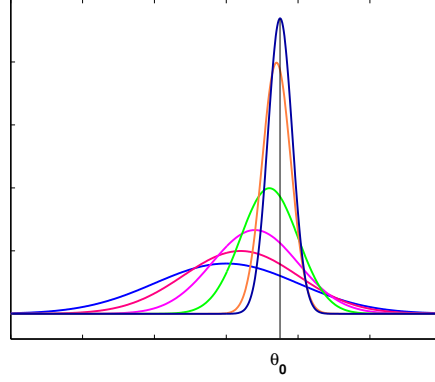
$$\text{cov}(\hat{\theta}_N) \leq \text{cov}(\bar{\theta}_N) \quad (4.4)$$

for all unbiased estimators  $\bar{\theta}_N$ . This means that it concerns an estimator that has the smallest possible variability (variance), measured in terms of its covariance matrix, of all unbiased estimators. Note that the inequality in (4.4) is a matrix inequality, requiring that the matrix  $\text{cov}(\bar{\theta}_N) - \text{cov}(\hat{\theta}_N)$  is positive semi-definite.

For scalar-valued estimators, the (relative) efficiency of an estimator  $\hat{\theta}_1$  with respect to another estimator  $\hat{\theta}_2$  is sometimes denoted by

$$\frac{\text{var}(\hat{\theta}_2)}{\text{var}(\hat{\theta}_1)}.$$





**Figure 4.1:** The most complete description of the properties of an estimator (which is a random variable), is its probability density function  $f_{\hat{\theta}_N}(\theta)$ , where  $\hat{\theta}$  is the estimate and  $N$  is the number of measurements used in the estimator.

---

**Example 4.1 (Estimating the variance of a random variable)**

Let  $\{\mathbf{x}_i\}_{i=1,\dots,N}$  be  $N$  independent observations of a random variable  $\mathbf{x}$ . We investigate the bias properties of the estimator

$$\widehat{\sigma^2} := \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2 \quad (4.5)$$

with  $\bar{\mathbf{x}} = (1/N) \sum_{i=1}^N \mathbf{x}_i$ , as an estimator of the mean of  $\mathbf{x}$ :

$$\sigma_{\mathbf{x}}^2 = E[(\mathbf{x} - E[\mathbf{x}])^2].$$

The mean value of the estimator can be analyzed as follows:

$$E[\widehat{\sigma^2}] = E\left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2\right] = \frac{1}{N} \sum_{i=1}^N E[(\mathbf{x}_i - \bar{\mathbf{x}})^2]. \quad (4.6)$$

Analyzing one term:

$$\begin{aligned} E[(\mathbf{x}_i - \bar{\mathbf{x}})^2] &= E\left[(\mathbf{x}_i - \frac{1}{N} \sum_{j=1}^N \mathbf{x}_j)^2\right] = E\left[(\mathbf{x}_i - \mu_{\mathbf{x}} + \mu_{\mathbf{x}} - \frac{1}{N} \sum_{j=1}^N \mathbf{x}_j)^2\right] \\ &= E\left[(\mathbf{x}_i - \mu_{\mathbf{x}} - \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \mu_{\mathbf{x}}))^2\right] \\ &= E[(\mathbf{x}_i - \mu_{\mathbf{x}})^2] + \frac{1}{N^2} \sum_{j=1}^N \sigma_{\mathbf{x}}^2 - 2E[(\mathbf{x}_i - \mu_{\mathbf{x}}) \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \mu_{\mathbf{x}})]. \end{aligned} \quad (4.7)$$

With respect to the last term on the right hand side of the last equation, we can write:

$$E[(\mathbf{x}_i - \mu_{\mathbf{x}}) \sum_{j=1}^N (\mathbf{x}_j - \mu_{\mathbf{x}})] = E[(\mathbf{x}_i - \mu_{\mathbf{x}})^2] + \underbrace{\sum_{j=1, j \neq i}^N E[(\mathbf{x}_i - \mu_{\mathbf{x}})(\mathbf{x}_j - \mu_{\mathbf{x}})]}_{=0},$$

where the second term is zero because of the fact that  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are independent random variable for  $i \neq j$ . As a result

$$E[(\mathbf{x}_i - \mu_{\mathbf{x}}) \sum_{j=1}^N (\mathbf{x}_j - \mu_{\mathbf{x}})] = \sigma_{\mathbf{x}}^2. \quad (4.8)$$

Substitution of this result in (4.7) delivers:

$$\begin{aligned} E[(\mathbf{x}_i - \bar{\mathbf{x}})^2] &= \sigma_{\mathbf{x}}^2 + \frac{1}{N^2} \cdot N\sigma_{\mathbf{x}}^2 - \frac{2}{N}\sigma_{\mathbf{x}}^2 \\ &= \sigma_{\mathbf{x}}^2(1 - \frac{1}{N}) = \frac{\sigma_{\mathbf{x}}^2(N-1)}{N}. \end{aligned} \quad (4.9)$$

As a result:

$$E[\widehat{\sigma^2}] = \frac{N-1}{N}\sigma_{\mathbf{x}}^2.$$

Consequently the estimator is biased for finite values of  $N$ . The bias disappears when  $N \rightarrow \infty$ , as in that case  $\frac{N}{N-1} \rightarrow 1$ . However the analysis also shows that one can construct an unbiased estimator

$$\widetilde{\sigma^2} := \frac{N}{N-1}\widehat{\sigma^2} = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2.$$

The division by  $N-1$  can be understood by realizing that one degree of freedom in the set of data is used to calculate  $\bar{\mathbf{x}}$ . As a result there remain  $N-1$  degrees of freedom to be used in the variance estimation.

---

## 4.2 Linear regression

### 4.2.1 Introduction

One of the most simple examples of estimation problems is the problem of estimating a linear relationship between two different (random) variables, on the basis of multiple observations of the two variables. Consider for instance the situation sketched in figure 4.2(left) where the (blue) dots reflect measurement pairs  $(u_i, y_i)$ . The question of finding a linear relationship between  $y$  and  $u$  can then be rephrased by finding the “best” straight line through the cloud of measurement points.

In order to solve this problem one generally has to make some assumptions on the source of the randomness in the data. We will first treat the standard case, where it is assumed that one variable ( $u$ ) is measured noise free, and the other variable ( $y$ ) is noise disturbed. Later we will comment on this and discuss generalizations of this paradigm. In the considered situation the relation between the measurements,  $\{u_i, y_i\}_{i=1, \dots, n}$ , is hypothesized by the model:

$$y_i = b_0 + b_1 u_i + e_i \quad (4.10)$$

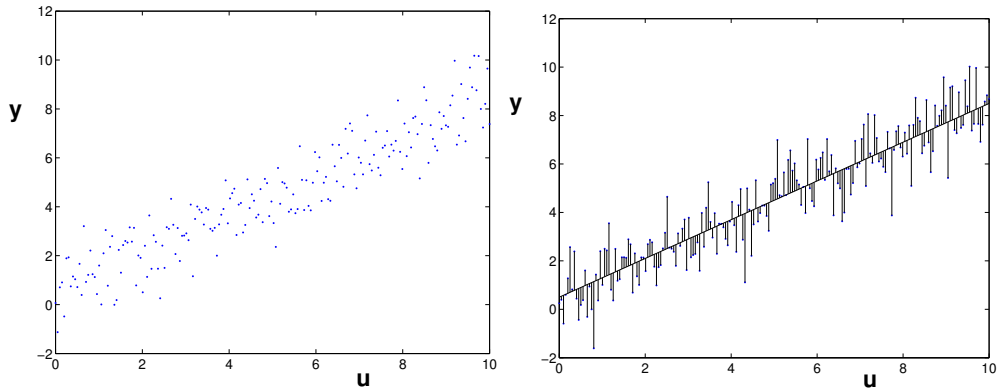
where  $b_0, b_1$  are unknown coefficients (parameters) that are to be estimated, and  $e_i$  is an error term that accounts for the fact that the measured points do not lie exactly on a straight line. The term  $e_i$  can be considered a realization of a random variable  $e$  with a particular probability density function. More attention will be paid to this when analyzing the properties of the linear least squares estimator and the weighted linear least squares estimator.

## 4.2.2 Linear least-squares estimation

The linear least-squares estimate, constructs a solution to the sketched problem, by looking for an estimate  $\hat{b}_0$  and  $\hat{b}_1$  such that

$$\sum_{i=1}^n e_i^2$$

is minimal. This implies that the “errors” between the straight line and measurement points are minimal (in squared sense) if one considers the “errors” to lie in the  $y$ -direction of the plot, as sketched in figure 4.2(right).



**Figure 4.2:** Observed data points  $\{u_i, y_i\}$  (left) and measure of fit to a straight line, by considering the errors in the  $y$  variables (right).

Starting from the model equation

$$y_i = b_0 + b_1 u_i + e_i \quad (4.11)$$

we write

$$y_i = \phi_i^T \theta + e_i \quad (4.12)$$

with

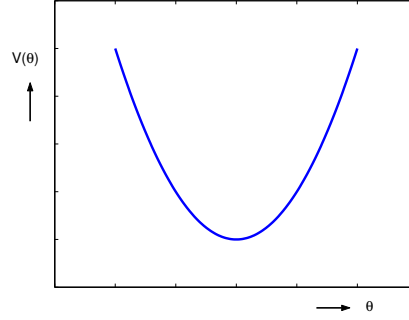
$$\phi_i = \begin{bmatrix} 1 \\ u_i \end{bmatrix} \text{ and } \theta = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}.$$

In this formulation the vector  $\phi_i$  is called the *regressor* or *regression variable*. The components of  $\phi_i$  are referred to as the *independent variables* in the regression problem, while the  $y_i$  are denoted the *dependent variables*.

The cost function  $J$  that has to be minimized is written as

$$J := \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \phi_i^T \theta)^2$$

When the data is given, the function  $J$  becomes a function of  $\theta$ ,  $J(\theta)$ . Besides the fact that  $J(\theta)$  is quadratic in  $e_i$ , it is also quadratic in  $\theta$ . This implies that it is a convex function, as illustrated in figure 4.3 for a scalar-valued  $\theta$ , having a unique global minimum that can be obtained by setting the derivative of  $J(\theta)$  to zero.



**Figure 4.3:**  $J$  is a quadratic function  $\theta$ , which can be visualized for the situation that  $\theta$  is scalar-valued.

Since<sup>3</sup>

$$\frac{\partial J(\theta)}{\partial \theta} = \begin{pmatrix} \frac{\partial J}{\partial b_0} \\ \frac{\partial J}{\partial b_1} \end{pmatrix} = \begin{pmatrix} -2 \sum_i (y_i - \phi_i^T \theta) \\ -2 \sum_i \phi_i (y_i - \phi_i^T \theta) \end{pmatrix}$$

it follows that

$$\left. \frac{\partial J}{\partial \theta} \right|_{\theta=\hat{\theta}} = 0 \rightarrow \sum_i \phi_i (y_i - \phi_i^T \hat{\theta}) = 0.$$

The resulting equations

$$\sum_i \phi_i (y_i - \phi_i^T \hat{\theta}) = 0 \quad (4.13)$$

are called the normal equations. Note that  $\phi_i$  and  $\theta$  are both 2-dimensional vectors, leading to a set of 2 (normal) equations with 2 unknowns. They deliver the following analytical solution for  $\hat{\theta}$ :

$$\left[ \sum_i \phi_i \phi_i^T \right] \hat{\theta} = \sum_i \phi_i y_i \quad (4.14)$$

or equivalently:

$$\hat{\theta} = \left[ \sum_i \phi_i \phi_i^T \right]^{-1} \sum_i \phi_i y_i \quad (4.15)$$

<sup>3</sup>As a notational convention, the derivative of a scalar with respect to a column vector is again a column vector.

provided that the inverse of the corresponding  $2 \times 2$ -matrix exists.

Note that the right hand side of the expression is only dependent on measurement data; once measurement data is available, the solution to the least-squares problem is simply obtained.

Existence of the matrix inverse in (4.15) is directly coupled to the question whether a sufficiently informative experiment has been done in order to uniquely determine the LS-solution. Consider, e.g., the situation that in the experiment all measurements  $u_i$  are the same, so

$$u_i = c \text{ for all } i.$$

Then  $\phi_i = (1 \ c)^T$  and

$$\sum_i \phi_i \phi_i^T = \sum_i \begin{bmatrix} 1 \\ c \end{bmatrix} [1 \ c] = \begin{pmatrix} n & cn \\ cn & nc^2 \end{pmatrix}$$

As the second column of this matrix is obtained by scalar multiplication of the first column by  $c$ , the matrix is singular and its inverse will not exist. Therefore the LS-solution will not be unique; there will exist many solutions of the equation (4.14). The non-uniqueness of the solution can also simply be understood by considering the problem in the scope of figure 4.2; if the cloud of points is concentrated around one value of  $u_i$ , then there does not exist a unique “best” straight line<sup>4</sup> that relates  $u$  to  $y$ .

### Least-squares solution in matrix form

For convenience, often use will be made of a more extensive matrix notation, replacing the summations of  $i$ .

The normal equations

$$\sum_i \phi_i (y_i - \phi_i^T \hat{\theta}) = 0 \quad (4.16)$$

can actually be rewritten as:

$$[\phi_1 \cdots \phi_n] \begin{pmatrix} y_1 - \phi_1^T \hat{\theta} \\ \vdots \\ y_n - \phi_n^T \hat{\theta} \end{pmatrix} = 0$$

which can be rewritten in short as

$$X^T (Y - X \hat{\theta}) = 0 \quad (4.17)$$

with

$$X = \begin{bmatrix} \phi_1^T \\ \vdots \\ \phi_n^T \end{bmatrix}; \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

---

<sup>4</sup>Note that the straight line  $u = c$  does not relate  $u$  to  $y$  and therefore is not a solution to the problem.

Note that  $X$  is a  $n \times 2$  matrix and  $Y$  an  $n \times 1$  vector. The solution is given by

$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

where again the assumption has to be made that the  $2 \times 2$  matrix inverse indeed exists.

### 4.2.3 Linear regression as a statistical estimation problem

Through imposing relation (4.11), the least-squares problem considered here has become a fully deterministic problem of drawing a straight line through a cloud of points, by minimizing the deviation to the line in the  $y$ -direction.

Considering the problem as a statistical estimation problem, we have to take care of the random character of the measurements  $y_i$ . If in (4.10) the variables  $e_i$  are random, then we have to write:

$$y_i = b_0 + b_1 u_i + e_i \quad (4.18)$$

showing that  $y_i$  actually are random variables as well.

The corresponding estimator of  $\theta_0$  then becomes

$$\hat{\theta} = (X^T X)^{-1} X^T \mathbf{Y},$$

with  $\mathbf{Y}$  now a vector of random variables:

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

and therefore  $\hat{\theta}$  also becomes a random variable. As a result we can study the statistical properties of the estimator  $\hat{\theta}$ .

### Bias of a linear least squares estimator

Assume that the measured data  $u_i, y_i$  is generated by an equation of the form

$$\mathbf{Y} = X\theta_0 + \mathbf{E} \quad \mathbf{E} = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}.$$

This implies that there exists a noise free output vector  $Y = X\theta_0$  that is observed through the random variable  $\mathbf{Y} = Y + \mathbf{E}$  where  $\mathbf{E}$  is a vector of random variables.

Substituting  $\mathbf{Y}$  into the expression for  $\hat{\theta}$ , it follows that

$$\begin{aligned} \hat{\theta} &= (X^T X)^{-1} X^T (X\theta_0 + \mathbf{E}) \\ &= \theta_0 + (X^T X)^{-1} X^T \mathbf{E}. \end{aligned} \quad (4.19)$$

The estimator is **unbiased** if  $E[\hat{\theta}] = \theta_0$ . If  $X$  is deterministic, as considered here, this condition is simply satisfied if  $E[\mathbf{E}] = 0$ , or equivalently if  $E[\mathbf{e}_i] = 0$  for all  $i$ .

Conclusion:

The linear least squares estimator is unbiased if the noise terms on the output variables are zero-mean random variables.

### Variance of a linear least squares estimator

Making use of the previous result, for an unbiased estimator we can write:

$$\text{cov}(\hat{\theta}) = E[(\hat{\theta} - \theta_0)(\hat{\theta} - \theta_0)^T] = E[(X^T X)^{-1} X^T \mathbf{E} \mathbf{E}^T X (X^T X)^{-1}].$$

If  $\mathbf{e}_i$  is a sequence of zero-mean uncorrelated random variables with equal variance  $\sigma^2$ , then

$$E[\mathbf{E} \mathbf{E}^T] = \sigma^2 \cdot I$$

and when the components of  $X$  are deterministic (no random variables), this will imply that

$$\text{cov}(\hat{\theta}) = (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} \quad (4.20)$$

$$= \sigma^2 (X^T X)^{-1}. \quad (4.21)$$

One of the important and appealing observations here is that the variance of  $\hat{\theta}$  increases with increasing values of  $\sigma^2$ : the higher the noise level, the larger the variance of the parameter estimates.

A graphical interpretation of the covariance matrix is obtained when considering the situation when  $\hat{\theta}$  has a Gaussian distribution with covariance matrix

$$\Sigma = \sigma^2 (X^T X)^{-1}.$$

In this case,

$$(\hat{\theta} - \theta_0)^T \Sigma^{-1} (\hat{\theta} - \theta_0) \in \chi_d^2,$$

which means that the left hand side expression follows a  $\chi^2$  distribution with  $d$  degrees of freedom (with  $d$  the dimension of  $\theta_0$ )<sup>5</sup>. The probability density function of a  $\chi_d^2$  distribution is

$$f(\hat{\theta}; d) = \frac{1}{2^{d/2} \Gamma(d/2)} x^{d/2-1} e^{-x/2} \quad (4.22)$$

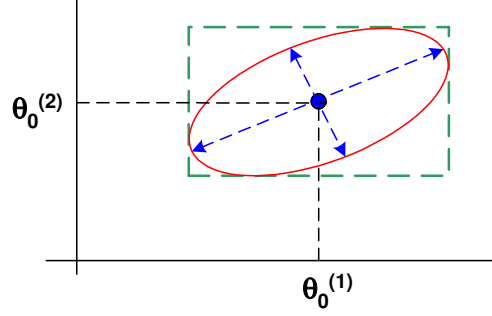
As a result the probability that

$$(\hat{\theta} - \theta_0)^T \Sigma^{-1} (\hat{\theta} - \theta_0) \geq \alpha$$

---

<sup>5</sup>A  $\chi_d^2$  distribution is typically obtained when taking the sum of  $d$  squared terms of independent normally distributed random variables with mean 0 and variance 1.

is specified by the  $\chi_d^2$  distribution. The sets determined by this latter expression determine ellipsoids in  $\mathbb{R}^2$ , as visualized in figure 4.4 (for the case  $d = 2$ ). Note that when giving confidence intervals for the two parameters on the basis of a particularly chosen level of probability  $\alpha$ , the covariance result (corresponding to the ellipsoidal area in figure 4.4) is less conservative than when bounding the two parameters separately (corresponding to the rectangular area in the figure).



**Figure 4.4:** Ellipsoid indicating levels of equal probability density function for a normally distributed estimator  $\hat{\theta}$  with covariance matrix  $\Sigma$ . The principal axes of the ellipsoid are determined by the eigenvectors and eigenvalues of  $\Sigma$ .

### Variance for incorrect models

The expression for the variance of  $\hat{\theta}$  remains the same if the chosen model does not match the underlying data generating equations. Consider, e.g., the situation where data originates from the relation

$$\mathbf{Y} = X\theta_0 + \mathbf{E}_0 = X_r\theta_0^{(1)} + X_e\theta_0^{(2)} + \mathbf{E}_0,$$

while a model is used with only a restricted number of regression variables:

$$\mathbf{Y} = X_r\theta + \mathbf{E}.$$

In this situation the regressor  $X_e$  is not incorporated in the model, possibly due to the fact that its relevance for the considered data was unknown to the user. Now the model estimator becomes:

$$\hat{\theta} = (X_r^T X_r)^{-1} X_r^T \mathbf{Y} = (X_r^T X_r)^{-1} X_r^T [X_r\theta_0^{(1)} + X_e\theta_0^{(2)} + \mathbf{E}_0] \quad (4.23)$$

$$= \theta_0^{(1)} + (X_r^T X_r)^{-1} X_r^T X_e \theta_0^{(2)} + (X_r^T X_r)^{-1} X_r^T \mathbf{E}_0. \quad (4.24)$$

If  $\mathbf{E}_0$  is zero-mean, and  $X_r$  and  $X_e$  are deterministic then

$$E[\hat{\theta}] = \theta_0^{(1)} + (X_r^T X_r)^{-1} X_r^T X_e \theta_0^{(2)},$$

and as a result

$$\text{cov}(\hat{\theta}) = E[(\hat{\theta} - E\hat{\theta})(\hat{\theta} - E\hat{\theta})^T] \quad (4.25)$$

$$= (X_r^T X_r)^{-1} X_r^T \sigma^2 I X_r (X_r^T X_r)^{-1} \quad (4.26)$$

$$= \sigma^2 (X_r^T X_r)^{-1}. \quad (4.27)$$



The resulting estimator will be biased now, but its variance expression still matches the general formula (4.21), where  $X$  then has to be interpreted as the regressor matrix that is used in the chosen model.

---

**Example 4.2 (Estimation of a physical variable from 5 different measurements)**

Consider a physical variable  $\theta_0$  that is measured by 5 different instruments, each having a different level of measurement noise. One can think, e.g., of the measurement of a temperature with different instruments. The available instruments are given by

$$\begin{aligned} \mathbf{y}_1 &= \theta_0 + \mathbf{e}_1 \\ \mathbf{y}_2 &= \theta_0 + \mathbf{e}_2 \\ \mathbf{y}_3 &= \theta_0 + \mathbf{e}_3 \\ \mathbf{y}_4 &= \theta_0 + \mathbf{e}_4 \\ \mathbf{y}_5 &= \theta_0 + \mathbf{e}_5 \end{aligned}$$

where  $\mathbf{e}_i$  is the random (additive) error that is induced by instrument number  $i$ . These “errors” are supposed to be zero-mean random variables with variance  $\sigma_i^2$ , such that

$$\sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_5.$$

There are several candidate estimators for estimating the true parameter value  $\theta_0$ :

- Choose the measurement of the “best” instrument, i.e., the instrument with the smallest variance error. This implies:

$$\hat{\theta} = \mathbf{y}_1. \quad (4.28)$$

This estimator is unbiased, and has a variance of  $\sigma_1^2$ .

- Combine all measurement into a simple linear least squares estimator. The corresponding model is:  $\mathbf{y}_i = \theta + \mathbf{e}_i$ , and the least squares solution is obtained by minimizing  $\sum_{i=1}^5 e_i^2$  is given by

$$\hat{\theta} = (X^T X)^{-1} X^T \mathbf{Y}$$

with  $X = [1 \dots 1]^T$  and  $\mathbf{Y} = [\mathbf{y}_1 \dots \mathbf{y}_5]^T$ , leading to

$$\hat{\theta} = \frac{1}{5} \sum_{i=1}^5 \mathbf{y}_i.$$

The estimator is simply obtained as the average of the 5 different instrument measurements. Note that again  $E\hat{\theta} = \theta_0$ , and so the estimator is again unbiased.

The variance of the estimator is determined by

$$\text{var}(\hat{\theta}) = E[(X^T X)^{-1} X^T \mathbf{E} \mathbf{E}^T X (X^T X)^{-1}]$$

which with  $X$  as given above, reduces to

$$\frac{1}{25} [1 \cdots 1] \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_5^2 \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

so that

$$\text{var}(\hat{\theta}) = \frac{1}{25} \sum_{i=1}^5 \sigma_i^2.$$

Note that it is not automatic that the variance of the second estimator is smaller than that of the first estimate. If all  $\sigma_i$  are equal, then averaging the 5 results improves the final variance with a factor 5; however when instruments 2 until 5 are much worse in quality than instrument number 1, averaging may even deteriorate the final estimator variance to become worse than the first one. If all  $\sigma_i$  are known, we can build an estimator that doesn't suffer from this deterioration in performance, see section 4.2.4.

#### Example 4.3 (Absorption coefficient of an optical fiber (Frieden, 2001))

A physical measurement problem does not always look as a regression problem, but can often be written in this format.

Consider the problem of experimentally determining the absorptance coefficient of an optical fiber. This can be done by measuring the light intensity  $i(x)$  at several measurement locations  $x$  along the fiber. The -theoretical- relation between light intensity and location is given by

$$I(x) = I_0 e^{-\alpha x}$$

with  $I(x)$  the light intensity at distance  $x$  from the light source,  $\alpha$  the absorption coefficient, and  $I_0$  the light intensity of the source.

If the initial intensity  $I_0$  is unknown, there are basically two unknown parameters:  $I_0$  and  $\alpha$ . Measurement of light intensity at location  $x$  is done by cutting the fiber, and measuring the light intensity. This is done for  $N$  different values of  $x$ .

In this situation it does not seem to be possible to write the system equation in the form:

$$I(x) = \phi^T \theta$$

as is required for a linear least squares estimator to be constructed.

However if we take the (natural) logarithm of the relation then:

$$\log I(x) = \log I_0 - \alpha x \quad (4.29)$$

$$= \phi^T \theta \quad (4.30)$$

with

$$\phi = \begin{pmatrix} 1 \\ -x \end{pmatrix} \quad \text{and} \quad \theta = \begin{pmatrix} \log I_0 \\ \alpha \end{pmatrix}.$$

By taking  $N$  measurements  $\{x_i, I(x_i)\}_{i=1, \dots, N}$  one can now construct a linear least squares estimator:

$$\hat{\theta}_N = \left[ \sum_i \phi_i \phi_i^T \right]^{-1} \sum_i \phi_i \log I(x_i).$$

In order for this estimator to be unbiased the noise on the measurement data has to be such that

$$\log I(x_i) = \log I_0 - \alpha x_i + e_i$$

with  $e_i$  realizations of a zero-mean stochastic process. In terms of the original relation for  $I(x)$ , this implies that there is a multiplicative noise contribution

$$I(x) = I_0 e^{-\alpha x} e^e$$

with  $e$  a zero-mean stochastic process.

#### 4.2.4 Weighted linear least squares estimation

In the linear regression problem the residual “error” at every measurement point is weighted with a constant weighting over all measurements. In other words: all data is equally weighted. In the weighted least squares approach, an additional weighting factor is introduced that allows for different weighting of the several measurements. This will be shown to provide means to arrive at estimators with improved statistical properties.

The *weighted least squares criterion* is formulated as:

$$J(\theta) = (Y - X\theta)^T W (Y - X\theta)$$

with  $W$  a symmetric positive definite  $n \times n$  matrix that is called the *weighting matrix*. The weighted least squares criterion may alternatively be written as

$$J(\theta) = \sum_i \sum_j w_{ij} (y_i - \phi_i^T \theta) (y_j - \phi_j^T \theta).$$

with  $\{w_{i,j}\}_{i,j=1, \dots, n}$  the elements of  $W$ . Note that for the special case of a diagonal weighting matrix  $W = \text{diag}(w_1, \dots, w_n)$ , the weighted least squares criterion can also be written as

$$J(\theta) = \sum_i w_i (y_i - \phi_i^T \theta)^2.$$

The weighted least squares solution is found by minimizing the weighted least squares criterion with respect to  $\theta$ . Setting the derivative of  $J(\theta)$  with respect to  $\theta$  to zero, the normal equations result:

$$X^T W (Y - X\hat{\theta}) = 0$$

If  $Y$  is again considered as a random variable (i.e.,  $Y = \mathbf{Y}$ ), the solution to the weighted least squares problem becomes:

$$\hat{\theta} = (X^T W X)^{-1} X^T W \mathbf{Y}$$

which reduces to the simple least squares estimator when  $W = I$ .

### Bias and variance of the weighted linear least squares estimator

For a data-generating system

$$\mathbf{Y} = X\theta_0 + \mathbf{E}_0$$

the weighted least squares estimator is given by:

$$\hat{\boldsymbol{\theta}} = (X^T W X)^{-1} X^T W \mathbf{Y} = \theta_0 + (X^T W X)^{-1} X^T W \mathbf{E}_0$$

so the estimator is unbiased if  $E[\mathbf{E}_0] = 0$ , i.e., if the noise process is zero-mean.

For the variance of the unbiased estimator one can write:

$$\begin{aligned} \text{cov}(\hat{\boldsymbol{\theta}}) &= E[(\hat{\boldsymbol{\theta}} - \theta_0)(\hat{\boldsymbol{\theta}} - \theta_0)^T] \\ &= E[(X^T W X)^{-1} X^T W \mathbf{E}_0 \mathbf{E}_0^T W X (X^T W X)^{-1}] \end{aligned}$$

Note that in the particular situation that  $W = (E[\mathbf{E}_0 \mathbf{E}_0^T])^{-1}$ , the expression for the variance simplifies to

$$\text{cov}(\hat{\boldsymbol{\theta}}) = (X^T W X)^{-1}. \quad (4.31)$$

This particular choice of weighting matrix is not just chosen for simplicity; it appears to have an important property. It is the - optimal - weighting that minimizes the variance over all possible choices of weighting matrices. More precisely, the weighted least squares estimator with as weighting matrix the inverse of the covariance matrix of the noise (i.e., with  $W = \Sigma^{-1}$ ) has minimum variance within the class of weighted least squares estimators, that is, the difference of the covariance matrix of any estimator of this class and its covariance matrix is positive semidefinite.

**Lemma 4.1.** *Let  $\Sigma$  be a positive definite matrix<sup>6</sup>, and define*

$$P(W) = (X^T W X)^{-1} X^T W \Sigma W X (X^T W X)^{-1}.$$

*Then for all symmetric, positive semi-definite  $W$ , it holds that*

$$P(\Sigma^{-1}) \leq P(W).$$

*This inequality expresses that the difference between the left-hand and right-hand member is negative semi-definite. A property of a negative semi-definite matrix is that its diagonal elements cannot be positive. This means that the diagonal elements of  $P(\Sigma^{-1})$  are smaller than or equal to the corresponding diagonal elements of  $P(W)$ .*

**Proof:** The matrix

$$\begin{bmatrix} X^T \\ X^T W \Sigma \end{bmatrix} \Sigma^{-1} \begin{bmatrix} X^T \\ X^T W \Sigma \end{bmatrix}^T = \begin{bmatrix} X^T \Sigma^{-1} X & X^T W X \\ X^T W X & X^T W \Sigma W X \end{bmatrix}$$

is positive semi-definite by construction.

---

<sup>6</sup>A matrix is positive definite if it's symmetric and all its eigenvalues are positive

Consider any positive semi-definite matrix

$$H = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}$$

with  $C$  invertible. Then considering  $xHx^T$  with  $x = [x_1 \quad -x_1BC^{-1}]$  and  $x_1$  arbitrary shows that

$$xHx^T = x_1[A - BC^{-1}B^T]x_1.$$

If  $H$  is positive semi-definite, then  $(A - BC^{-1}B^T)$  is also positive semi-definite. Applying this result to the above expressions, shows that

$$X^T\Sigma^{-1}X - X^TWX[X^TW\Sigma WX]^{-1}X^TWX \geq 0$$

or equivalently

$$X^T\Sigma^{-1}X \geq X^TWX[X^TW\Sigma WX]^{-1}X^TWX.$$

Taking inverses of both sides of the inequality leads to

$$(X^T\Sigma^{-1}X)^{-1} \leq (X^TWX)^{-1}[X^TW\Sigma WX](X^TWX)^{-1}.$$

This reduces to equality when  $W = \Sigma^{-1}$  i.e. the smallest possible covariance matrix.

□

A weighted least-squares estimator, with the optimal weighting as indicated above, is also referred to as the [Markov estimator](#) or the [Best Linear Unbiased Estimator \(BLUE\)](#). The reason for this is that it can be shown that this estimator has the minimum variance not only within the class of weighted least squares estimators, but also within the (broader) class of estimators that are linear in the observations ( $\mathbf{Y}$ ) and are unbiased.

If we apply this BLUE to the problem as sketched in Example 4.2, then the “optimal” weighting should be chosen as

$$W = \Sigma^{-1} = \begin{bmatrix} \sigma_1^{-2} & & \\ & \ddots & \\ & & \sigma_5^{-2} \end{bmatrix}$$

The weighted least squares criterion now reads:

$$\min_{\theta} \sum_{i=1}^5 \frac{(y_i - \theta)^2}{\sigma_i^2},$$

leading to the weighted least squares solution

$$\begin{aligned} \hat{\theta} &= (X^T\Sigma^{-1}X)^{-1}X^T\Sigma^{-1}\mathbf{Y} \\ &= \frac{[y_1/\sigma_1^2 + \dots + y_5/\sigma_5^2]}{1/\sigma_1^2 + \dots + 1/\sigma_5^2}. \end{aligned} \tag{4.32}$$

The variance of this estimator is given by (4.31):

$$\begin{aligned} \text{var}(\hat{\theta}) &= (X^T \Sigma^{-1} X)^{-1} \\ &= \frac{1}{1/\sigma_1^2 + \dots + 1/\sigma_5^2}. \end{aligned} \quad (4.33)$$

Consequently: the use of all five measurement devices reduces the variance of the final temperature estimator, and is to be preferred over simply choosing the measurement from the device with the highest precision. However in order to find the optimal weightings that minimize the overall variance, knowledge of the variance of the several instruments is required.

---

### 4.3 The Cramér-Rao lower bound

The estimators presented so far show different properties for the variance (i.e., the covariance matrix) of the estimated parameter. This raises the fundamental question whether, in a particular estimation problem, there exists a lower bound on the reachable estimator variance. The answer to this question is one of the most powerful results in estimation theory, and is known as the Cramér-Rao lower bound (CRLB), named after the work of Cramér (1946) and Rao (1945)<sup>7</sup>.

#### Cramér-Rao lower bound (CRLB)

Consider observations from a random variable  $\mathbf{y}$  with probability density function  $f_{\mathbf{y}}(\mathbf{y}, \theta)$ , where  $\theta$  is the unknown parameter. Then for *any* unbiased estimator  $\hat{\theta}$  of the parameter  $\theta$ , its covariance matrix satisfies the inequality<sup>8</sup>

$$\text{cov}(\hat{\theta}) \geq F^{-1} \quad (4.34)$$

with the *Fisher Information Matrix*:

$$F = E \left[ - \frac{\partial^2}{\partial \theta^2} \log f_{\mathbf{y}}(\mathbf{y}; \theta) \Big|_{\theta=\theta_0} \right] \quad (4.35)$$

The logarithm of the probability density function of the measurement data determines the lower bound on the variance of any unbiased parameter estimator.

Note that the measurable random variable  $\mathbf{y}$  will generally be an  $n$ -dimensional vector, and the corresponding probability density function  $f_{\mathbf{y}}$  a multivariate

<sup>7</sup>Actually Fisher (1922) had shown a similar result two decades earlier.

Ronald Aylmer Fisher (1890-1962) was a British mathematician who played a key role in the development of modern probability theory.

<sup>8</sup>The expression for the Fisher Information Matrix requires the second partial derivative of  $\log f$  (the natural logarithm) with respect to  $\theta$  to exist and to be absolutely integrable.

p.d.f. With  $\theta$  being a  $d$ -dimensional vector of parameters, the second partial derivative of a scalar function  $g$  with respect to  $\theta$  is a  $d \times d$  matrix defined by

$$\frac{\partial^2}{\partial \theta^2} g = \begin{pmatrix} \frac{\partial^2 g}{\partial \theta_1 \partial \theta_1} & \frac{\partial^2 g}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 g}{\partial \theta_1 \partial \theta_d} \\ \frac{\partial^2 g}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 g}{\partial \theta_2 \partial \theta_2} & \cdots & \frac{\partial^2 g}{\partial \theta_2 \partial \theta_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 g}{\partial \theta_d \partial \theta_1} & \frac{\partial^2 g}{\partial \theta_d \partial \theta_2} & \cdots & \frac{\partial^2 g}{\partial \theta_d \partial \theta_d} \end{pmatrix}.$$

### Proof of the Cramér-Rao inequality

Since  $f_{\mathbf{y}}(y; \theta)$  is a probability density function, it holds that

$$\int_{-\infty}^{\infty} f_{\mathbf{y}}(y; \theta) dy = 1 \quad \text{for all } \theta \quad (4.36)$$

where the integral over  $dy$  should be considered to be taken over the  $n$ -dimensional measurement space, when the measurement data is  $n$ -dimensional. In particular it will follow that

$$\int_{-\infty}^{\infty} \theta f_{\mathbf{y}}(y; \theta) dy = \theta \quad \text{for all } \theta. \quad (4.37)$$

The expected value of an estimator  $\hat{\theta}(\mathbf{y})$  is given by

$$E[\hat{\theta}(\mathbf{y})] = \int_{-\infty}^{\infty} \hat{\theta}(y) f_{\mathbf{y}}(y; \theta_0) dy.$$

Consequently one can write that for every unbiased estimator  $\hat{\theta} = \hat{\theta}(\mathbf{y})$  it follows that

$$E[\hat{\theta} - \theta_0] = \int_{-\infty}^{\infty} [\hat{\theta}(y) - \theta_0] f_{\mathbf{y}}(y; \theta_0) dy = 0.$$

Basically this expression holds true for all  $\theta_0$ , as the unbiasedness of the estimator  $\hat{\theta}$  is not dependent on the particular choice of  $\theta_0$ . As a result the above equation should also hold when differentiated with respect to  $\theta_0$ .

If the partial derivative  $\frac{\partial f_{\mathbf{y}}(y; \theta)}{\partial \theta}$  exists and is absolutely integrable, then

$$\frac{\partial}{\partial \theta_0^T} \int_{-\infty}^{\infty} [\hat{\theta}(y) - \theta_0] f_{\mathbf{y}}(y; \theta_0) dy = \quad (4.38)$$

$$= \int_{-\infty}^{\infty} [\hat{\theta}(y) - \theta_0] \left[ \frac{\partial f_{\mathbf{y}}(y; \theta_0)}{\partial \theta_0} \right]^T dy - \int_{-\infty}^{\infty} I \cdot f_{\mathbf{y}}(y; \theta_0) dy = 0 \quad (4.39)$$

showing that

$$\int_{-\infty}^{\infty} [\hat{\theta}(y) - \theta_0] \left[ \frac{\partial f_{\mathbf{y}}(y; \theta_0)}{\partial \theta_0} \right]^T dy = I.$$

With the expression for the derivative of the logarithm, this leads to

$$\int_{-\infty}^{\infty} [\hat{\theta}(y) - \theta_0] \left[ \frac{\partial \log f_{\mathbf{y}}(y; \theta_0)}{\partial \theta_0} \right]^T f_{\mathbf{y}}(y; \theta_0) dy = I,$$

or

$$E \left[ [\hat{\boldsymbol{\theta}} - \theta_0] \frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta^T} \Big|_{\theta=\theta_0} \right] = I. \quad (4.40)$$

In the scalar case ( $d = 1$ ) one can now take the square of this relation, and apply Schwartz inequality  $(E[\mathbf{x}\mathbf{y}])^2 \leq E[\mathbf{x}^2] \cdot E[\mathbf{y}^2]$ , to show that

$$1 \leq E[(\hat{\boldsymbol{\theta}} - \theta_0)^2] \cdot E \left[ \left( \frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta) \Big|_{\theta=\theta_0} \right)^2 \right]$$

leading to

$$E[\hat{\boldsymbol{\theta}} - \theta_0]^2 \geq \left\{ E \left[ \left( \frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta) \Big|_{\theta=\theta_0} \right)^2 \right] \right\}^{-1}. \quad (4.41)$$

In the multivariate situation ( $d > 1$ ) a different route has to be followed. In that situation we rewrite (4.40) into the form

$$E[\mathbf{a}\mathbf{b}^T] = I \quad (4.42)$$

and we consider

$$E \left[ \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}^T \right] = \begin{pmatrix} E[\mathbf{a}\mathbf{a}^T] & I \\ I & E[\mathbf{b}\mathbf{b}^T] \end{pmatrix}$$

which is a positive semi-definite matrix by construction. Then with the property of positive semi-definite matrices, as explained in the proof of Lemma 4.1 it follows that

$$E[\mathbf{a}\mathbf{a}^T] \geq (E[\mathbf{b}\mathbf{b}^T])^{-1},$$

leading to

$$\text{cov}(\hat{\boldsymbol{\theta}}) \geq \left\{ E \left[ \left( \frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta) \Big|_{\theta=\theta_0} \right) \left( \frac{\partial}{\partial \theta} \log f_{\mathbf{y}}(\mathbf{y}; \theta) \Big|_{\theta=\theta_0} \right)^T \right] \right\}^{-1}. \quad (4.43)$$

Differentiating (4.36) shows that

$$\int_{-\infty}^{\infty} \frac{\partial f_{\mathbf{y}}(y; \theta)}{\partial \theta} dy = 0$$

or equivalently

$$\int_{-\infty}^{\infty} \frac{\partial \log f_{\mathbf{y}}(y; \theta)}{\partial \theta} f_{\mathbf{y}}(y; \theta) dy = 0.$$

Differentiating this equation with respect to  $\theta^T$  gives

$$\int_{-\infty}^{\infty} \left[ \frac{\partial^2 \log f_{\mathbf{y}}(y; \theta)}{\partial \theta^2} + \left( \frac{\partial \log f_{\mathbf{y}}(y; \theta)}{\partial \theta} \right) \left( \frac{\partial \log f_{\mathbf{y}}(y; \theta)}{\partial \theta} \right)^T \right] f_{\mathbf{y}}(y; \theta) dy$$



and as a result

$$E \left[ \left( \frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta} \right) \left( \frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta} \right)^T \right] = -E \left[ \frac{\partial^2 \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta^2} \right].$$

Together with (4.43) this proves the result for the Cramér-Rao bound.  $\square$

### Discussion

Some comments on the nature of the lower bounds:

- The calculation of the CRLB requires exact knowledge of the probability density function of the measured random variables. This is of course a quite strict requirement, that asks for detailed knowledge of the disturbances that act on the measured data.
- The calculation of the CRLB will generally require knowledge of the exact system parameter  $\theta_0$ . There are exceptions though, when the second derivative of the  $\log f_{\mathbf{y}}(\mathbf{y}; \theta)$  is independent of  $\theta$ . This typically happens in the situation of Gaussian probability density functions, in combination with the requirement that the parameter  $\theta$  occurs linearly in the measured output. See also Example 4.4.
- The CRLB is a lower bound on the variance of unbiased estimators. It does not say anything about the variance of biased estimators, and so it is possible that there exists a biased estimator that has a smaller variance, and possibly even a smaller mean-squared error.

As a result of these remarks, calculation of the CRLB is often not feasible in practical situations. However the bound is very useful in analysis questions, e.g., when evaluating properties of several estimators and when comparing the covariance matrices of several estimators.

Notice that the CRLB is not related to a particular estimation method. It depends on the statistical properties of the observed variables, the measured data, and in most cases the hypothetical true values of the parameters. As first sight, this dependence on the true values looks as a serious impediment to the practical use of the bound. However, the expressions for the bound provide the experimenter with the means to compute numerical values for it, using nominal values of the parameters. This provides the experimenter with quantitative insight in what precision (i.e., variance) might be achieved from the available observations. In addition, it provides insight in the sensitivity of the precision to the parameter values.

Another important purpose for which the expressions for the CRLB can be used, is the optimization of the experiment design, i.e., the selection of which data to measure and to use as a basis for the estimation. By calculating the CRLB, the experimenter gets an impression if for a given experiment setup the precision attainable is sufficient for the purpose concerned. If not, the experiment design has to be changed. If this is not possible, it has to be concluded that the observations are not suitable for the purpose of the measurement procedure. In this way,

the experiment design can be optimized so as to attain the highest precision, i.e., the smallest variance (see, e.g., Van den Bos, 1999).

As a final remark it should be noted that existence of a lower bound on the parameter variance, does not imply that an estimator can be found that reaches this lower bound. Especially in the situation of observing random variables with a finite number of observations, it appears very hard to find the minimum variance estimator.

---

**Example 4.4 (Continuation of Example 4.2)**

Consider again the 5 measurements

$$\mathbf{y}_i = \theta_0 + \mathbf{e}_i \quad i = 1, \dots, 5.$$

For calculation of the CRLB of this estimation problem, we additionally assume that the random errors  $\mathbf{e}_i$  are jointly Gaussian distributed, with mean value 0 and covariance matrix  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_5^2)$ .

Using the multivariate Gaussian distribution it follows that

$$f_{\mathbf{y}}(\mathbf{y}, \theta) = \frac{1}{(2\pi)^{5/2} \sqrt{\det \Sigma}} \exp\left[-\frac{1}{2}(\mathbf{y} - \underline{\theta})^T \Sigma^{-1}(\mathbf{y} - \underline{\theta})\right]$$

with  $\mathbf{y} := [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_5]^T$ , and  $\underline{\theta} = [\theta \ \dots \ \theta]^T$ . Consequently

$$\log f_{\mathbf{y}}(\mathbf{y}, \theta) = c - \sum_{i=1}^5 \frac{1}{2} \frac{(\mathbf{y}_i - \theta)^2}{\sigma_i^2}.$$

Taking the partial derivative with respect to  $\theta$  delivers

$$\frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta} = \sum_{i=1}^5 \frac{(\mathbf{y}_i - \theta)}{\sigma_i^2}$$

and the second derivative:

$$\frac{\partial^2 \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta^2} = \sum_{i=1}^5 \frac{-1}{\sigma_i^2}.$$

Substituting this result in the expression for the CRLB, gives:

$$\text{var}(\hat{\theta}) \geq \frac{1}{1/\sigma_1^2 + \dots + 1/\sigma_5^2}.$$

This provides us, in the considered situation, with the best possible variance for any unbiased estimator.

Note that the weighted least squares estimator that was analyzed in section 4.2.4 reaches exactly this lower bound of the parameter variance in expression (4.33).

The conclusions that one can draw from this, are the following

- In the problem setting of the considered example, the weighted least squares estimator with the particular weights as chosen in section 4.2.4 leads to the smallest possible variance among all unbiased weighted least squares estimators, irrespective of the probability density function of the disturbances.
- If in the problem setting of the considered example, the disturbance terms are Gaussian distributed, then the weighted least squares estimator with the above mentioned weights, has the smallest possible variance among *all* unbiased estimators.

If the probability density function of the disturbances is unknown, the weighted least squares estimator can lead to satisfactory results; however there is no guarantee that it is the best possible estimator in terms of minimum variance.

---

## 4.4 Maximum likelihood estimator

In the previous section it was shown that the probability density function of the measured random variables plays a crucial role in the smallest possible variance that can be reached for any unbiased parameter estimator.

The parameter estimation methods considered so far are based on linear (regression) techniques, and therefore they are restricted to a well-defined class, characterized by the fact that the resulting estimators lead to simple analytical functions of the measured random variables. As a result, the computational tools required for the estimators are very straightforward.

In this section a general philosophy to parameter estimation will be presented that takes account of the probability density function of the measured random variables, and therefore has the potentials to come close to the covariance expressions as induced by the CRLB.

The so-called maximum likelihood principle is based on the following reasoning.

Suppose that a (vector) random variable  $\mathbf{y}$  is observed, and that the underlying probability density function of  $\mathbf{y}$  is given by

$$f_{\mathbf{y}}(y, \theta)$$

where  $\theta$  reflects the unknown (vector) parameter that is to be estimated.

For a given value of  $\theta$  the function  $f_{\mathbf{y}}(y, \theta)$  is a probability density function. However for a fixed value of  $y$  and unknown  $\theta$ , the function  $f$  is a (deterministic) function of  $\theta$  and referred to as the *likelihood function*, indicated by

$$L(\theta; y).$$

### Maximum likelihood principle

For a given observation  $y$  of the measured random variable  $\mathbf{y}$ , determine the maximum likelihood estimate as that value of  $\theta$  that maximizes the likelihood function, i.e.,

$$\max_{\theta} L(\theta; y).$$

The resulting maximum likelihood (ML) estimator is denoted as:

$$\hat{\theta}_{ml} = \arg \max_{\theta} L(\theta; \mathbf{y}).$$

For given observed values  $\mathbf{y} = y$ , the maximum likelihood estimate is determined by that value of  $\theta$  for which the probability density function  $f_{\mathbf{y}}(y, \theta)$  reaches its maximum value. In other words: that parameter  $\theta$  is sought for that generates a probability density function for which the observed measurement data was the most probable data. The idea can best be visualized in a very simple example.

---

#### Example 4.5 (Visualization of the maximum likelihood principle)

We consider the linear model between the observed variables  $\mathbf{y}$  and  $u$ , where  $u$  is exactly measured and known, and  $\mathbf{y}$  is a random variable that satisfies :

$$\mathbf{y} = \theta_0 u + \mathbf{e}$$

with  $\mathbf{e}$  a random variable with a particular pdf  $f_{\mathbf{e}}$ , and  $\theta_0$  an unknown scalar constant.

We are going to estimate a model of the form

$$\mathbf{y} = \theta u + \mathbf{e}$$

and the pdf of the observed random variable  $\mathbf{y}$  is then given by  $f_{\mathbf{e}}(y - \theta u)$ .

For one observed measurement of  $\mathbf{y}$ , i.e.,  $\mathbf{y}$  is a scalar, the pdf of  $\mathbf{y}$  as a function of  $\theta$  is depicted in Figure 4.5 for  $u = 2$ . It is a continuum of pdf's, since  $\theta$  varies over a continuous region.

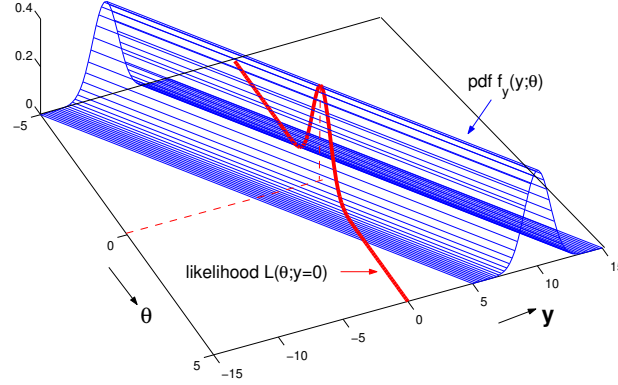
Now, if one observation of  $\mathbf{y}$  is made, e.g.,  $y = 0$ , then the likelihood function results:  $L(\theta) = f_{\mathbf{e}}(y - \theta u)|_{y=0}$ , being a function of  $\theta$  only, and depicted in Figure 4.5 as the solid (red) curve. The maximum likelihood estimate of  $\theta$  is that value of  $\theta$  for which  $L$  reaches its maximum.

In Figure 4.5 the example is sketched of a normal pdf  $f_{\mathbf{e}}$  with zero mean and unit variance. The likelihood function then becomes

$$L(\theta; y) = \frac{1}{\sqrt{2\pi}} e^{\frac{-(y-2\theta)^2}{2}}.$$

For the observation  $y = 0$ , maximization of  $L(\theta; y = 0)$  comes down to maximizing

$$\frac{1}{\sqrt{2\pi}} e^{\frac{-4\theta^2}{2}}$$



**Figure 4.5:** A prior probability density function of  $y$ , as a function of  $\theta$ , (blue continuum of curves), and the likelihood function  $L$  for observation  $y = 0$  (red solid curve).

which is obtained for  $\theta = 0$ , this being the maximum likelihood estimate.

### Maximum likelihood estimation of the parameters of linear regression models from Gaussian distributed observations

The considered -very simple- example can straightforwardly be extended to the situation of multiple observations and vector parameters. Consider the previously used linear regression model:

$$y_i = \phi_i^T \theta + e_i \quad i = 1, \dots, n$$

with

$$\phi_i = \begin{bmatrix} 1 \\ u_i \end{bmatrix} \quad \text{and} \quad \theta = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}.$$

and assume that  $e_i$  is a set of independent Gaussian random variables with pdf  $f_e$  having mean value 0 and variance  $\sigma^2$ .

Then the joint probability density function of the measured variables  $y_1 \cdots y_n$  is given by

$$f_y(y) = \prod_{i=1}^n f_e(y_i - \phi_i^T \theta)$$

and the likelihood function then becomes

$$L(\theta; Y) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \phi_i^T \theta)^2}{2\sigma^2}}.$$

Maximizing the likelihood function over  $\theta$  leads to the same argument as maximizing  $\log L(\theta, Y)$  over  $\theta$ . This is due to the fact that the log is a monotone increasing function. Taking the logarithm is often advantageous, in particular for

probability density functions that contain exponentials. Instead of maximizing  $\log L$  one can equivalently minimize  $-\log L$ , which is given by

$$-\log L(\theta; Y) = \frac{n}{2} \log 2\pi + n \log \sigma + \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \phi_i^T \theta)^2. \quad (4.44)$$

Since the first two terms on the right hand side of this expression are constants and not functions of  $\theta$ , consequently

$$\hat{\theta}_n = \arg \min_{\theta} \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \phi_i^T \theta)^2 \quad (4.45)$$

$$= \arg \min_{\theta} \sum_{i=1}^n (\mathbf{y}_i - \phi_i^T \theta)^2 \quad (4.46)$$

and this expression is exactly the same as the simple least squares (linear regression) estimator that was developed in section 4.2. This leads to the following conclusion.

For  $n$  independent observations from a Gaussian distribution with equal variance for all observations, the ML estimator is given by the simple least squares (LS) estimator.

If in the considered problem the noise terms  $\mathbf{e}_i$  are independent Gaussian random variables with zero mean and with fixed and known variance  $\sigma_i^2$ , being different for the different measurements  $i$ , then the likelihood function simply generalizes to

$$L(\theta; Y) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(\mathbf{y}_i - \phi_i^T \theta)^2}{2\sigma_i^2}}$$

and the corresponding maximum likelihood estimator is given by

$$\hat{\theta}_n = \arg \min_{\theta} \sum_{i=1}^n \frac{(\mathbf{y}_i - \phi_i^T \theta)^2}{\sigma_i^2} \quad (4.47)$$

and this expression is the same as the weighted least squares estimator with a diagonal weighting matrix  $W = \text{diag}(1/\sigma_1^2, \dots, 1/\sigma_n^2)$ , which was developed in section 4.2.4. If in the considered problem the noise terms  $\mathbf{e}_i$  are *correlated*

Gaussian random variables with zero mean and  $n \times n$  covariance matrix  $\Sigma$ , then the likelihood function generalizes to

$$L(\theta; Y) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\Sigma)}} e^{-\frac{1}{2}(Y - X\theta)^T \Sigma^{-1} (Y - X\theta)}$$

with

$$X = \begin{bmatrix} \phi_1^T \\ \vdots \\ \phi_n^T \end{bmatrix}; \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

and the corresponding maximum likelihood estimator is given by

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta}} (\mathbf{Y} - X\boldsymbol{\theta})^T \Sigma^{-1} (\mathbf{Y} - X\boldsymbol{\theta}), \quad (4.48)$$

which leads to

$$\hat{\boldsymbol{\theta}}_n = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \mathbf{Y} \quad (4.49)$$

and this expression is the same as the weighted least squares estimator that was developed in section 4.2.4. As a result one can state the following:

For  $n$  observations from a joint Gaussian distribution, the ML estimator is given by the weighted least squares (WLS) estimator, where the weighting matrix is given by the inverse of the covariance matrix of the observations.

Note that the considered weighted least squares estimator is identical to the BLUE (best linear unbiased estimator), discussed in section 4.2.4. So, for Gaussian disturbances (and regression models linear in the unknown parameters  $\boldsymbol{\theta}$ ), the BLUE and the ML estimator coincide, and from the results of the previous section it follows that their variance reaches the Cramér-Rao lower bound, i.e., the minimum possible variance over all unbiased estimators.

Note that the equivalence between ML and LS/WLS estimators typically holds true for Gaussian distribution functions. However the maximum likelihood principle goes beyond this situation and applies also to other distributions.

### Discussion

In the general case, the ML estimator will require solving the optimization problem

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta}} (-\log L(\boldsymbol{\theta}, \mathbf{Y}))$$

which might not be straightforward at all. In particular for Gaussian distributions, and a linear regression model, this optimization problem is simple, since it reduces to minimizing a quadratic function in  $\boldsymbol{\theta}$ . This convex optimization problem is simply solvable by efficient algorithms, relying on analytical expressions for its solution. In the general situation however the optimization problem will require nonlinear (gradient-type) optimization tools that suffer from possible poor convergence due to the existence of local minima, and the resulting lack of guarantee that a solution to the problem has been obtained. For an overview of optimization methods see, e.g., Fletcher (1980) and Miller (2000).

### Properties of the ML estimator

The maximum likelihood estimator has several important properties:

Under general conditions, the ML-estimator has the property that for the number of observations  $n$  tending to infinity,

$$\hat{\theta}_n \rightarrow \mathcal{N}(\theta_0, F^{-1})$$

meaning that the random variable  $\hat{\theta}_n$  converges in distribution to a Gaussian distribution with mean value  $\theta_0$  and covariance matrix equal to the Cramér-Rao bound (see, e.g., Van den Bos, 2007).

As a result of this, the ML-estimator is

- asymptotically unbiased;
- consistent
- asymptotically efficient, i.e., it asymptotically reaches the minimum possible variance (CRLB) among all unbiased estimators.

These properties of the ML estimator are very powerful, and are a strong support for this estimation principle. With respect to the practical implications one should realize though that in order for an ML estimator to be applicable, detailed knowledge of the probability density function of the observations has to be available (as prior information, e.g., derived for physical considerations/analysis of the disturbances that are acting on the measurement data). Note also that the attractive properties of the ML estimator all hold asymptotically in the number of data  $n$ . There is no guarantee for unbiased and minimum variance estimators in the case of finite data.

This section on the ML estimator will be ended by providing some examples.

---

#### Example 4.6 (Estimation of mean value of Gaussian random variable)

Consider  $n$  independent observations of a Gaussian random variable  $\mathbf{y}$  with unknown mean  $\mu$  and known variance  $\sigma^2$ . The likelihood function, as a function of  $\mu$  is given by

$$L(\mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mathbf{y}_i - \mu)^2}{2\sigma^2}} \quad (4.50)$$

and

$$\log L(\mu) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \mu)^2. \quad (4.51)$$

Maximizing  $L$  leads to the same argument as minimizing  $-\log L$ , and so the



maximum likelihood estimator for  $\mu$  is determined by

$$\hat{\mu}_n = \arg \min_{\mu} \left\{ -n \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \mu)^2 \right\} \quad (4.52)$$

$$= \arg \min_{\mu} \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \mu)^2. \quad (4.53)$$

Setting the derivative of the function in the right hand side to zero:

$$\left[ -2 \sum_{i=1}^n (\mathbf{y}_i - \mu) \right]_{\mu=\hat{\mu}_n} = 0 \quad (4.54)$$

delivers  $n\hat{\mu}_n = \sum_{i=1}^n \mathbf{y}_i$  or equivalently

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i. \quad (4.55)$$

In other words, in the considered situation the sample average is the maximum likelihood estimator for the mean value of the random variable.

---

#### Example 4.7 (Estimation of variance of Gaussian random variable)

Consider  $n$  independent observations of a Gaussian random variable  $\mathbf{y}$  with known mean  $\mu$  and unknown variance  $\sigma^2$ . The likelihood function, as a function of  $\sigma$  is given by

$$L(\sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mathbf{y}_i - \mu)^2}{2\sigma^2}}.$$

Since maximizing  $L(\sigma)$  has the same optimum argument as minimizing  $-\log L(\sigma)$ , we consider

$$-\log L(\sigma) = \frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{y}_i - \mu)^2 - n \log \frac{1}{\sqrt{2\pi}\sigma}$$

Minimizing  $-\log L(\sigma)$  can be done by setting  $\frac{\partial}{\partial \sigma} = 0$ , i.e.,

$$\left[ \frac{n}{\sigma} - 2 \sum_{i=1}^n \frac{(\mathbf{y}_i - \mu)^2}{2\sigma^3} \right]_{\sigma=\hat{\sigma}_n} = 0$$

or equivalently

$$n - \sum_{i=1}^n \frac{(\mathbf{y}_i - \mu)^2}{\hat{\sigma}_n^2} = 0$$

from which follows that

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \mu)^2. \quad (4.56)$$

The difference of this estimator with analyzed in Example 4.1 is that here the exact expression  $\mu$  is used, whereas in the earlier example an estimate of  $\mu$  was substituted in the expression. The current estimator is unbiased, whereas the estimator in (4.5) is biased.

---

#### Example 4.8 (Estimating arrival rate of photons on a detection plate)

The time  $T$  between the detection of two photons on a detection plate of an electron microscope is a random variable with exponential density function

$$f_{\mathbf{T}}(T) = \begin{cases} \alpha e^{-\alpha T} & T \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (4.57)$$

The coefficient  $\alpha$  is known as the arrival rate. Additionally it is assumed that intervals between arrivals are independent. For the estimation of the arrival rate  $\alpha$ , we construct the likelihood function:

$$L(\alpha; T_1, \dots, T_n) = \prod_{i=1}^n \alpha e^{-\alpha T_i} = \alpha^n e^{-\alpha \sum_{i=1}^n T_i}. \quad (4.58)$$

The maximum likelihood estimate is found by taking the logarithm of this function, and setting the derivative equal to 0:

$$\frac{\partial}{\partial \alpha} \left[ n \log \alpha - \alpha \sum_{i=1}^n T_i \right] = \frac{n}{\alpha} - \sum_{i=1}^n T_i = 0. \quad (4.59)$$

Solving this equation for  $\alpha$  leads to the maximum likelihood estimate

$$\hat{\alpha}_{ml} = \frac{1}{\frac{1}{n} \sum_{i=1}^n T_i}. \quad (4.60)$$

The estimate is the inverse of the average interarrival time.

---

In the examples given so far, the maximum likelihood estimate can always be obtained by simple analytic expressions. This is however not a structural property of this estimator.

Consider, e.g., the emission of radioactive particles from two radioactive sources, generally modelled as

$$y_t = \lambda_1 e^{\mu_1 t} + \lambda_2 e^{\mu_2 t} + \eta + e_t$$

with  $\lambda_1, \lambda_2$  the concentration of the two radioactive sources,  $\mu_1, \mu_2$  the decay rate of the separate sources,  $\eta$  the effect of background radiation, and  $e_t$  a random error, and that the radiation  $y_t$  is observed over  $n$  time instants with  $t = n \cdot \delta t$ .

Suppose that  $e_t$  are independent Gaussian random variables with known and fixed mean and variance.

If  $\mu_1$  and  $\mu_2$  are known, and we intend to estimate  $\lambda_1, \lambda_2$  and  $\eta$ , we can write

$$y_t = \phi_t^T \theta + e_t \quad (4.61)$$

with

$$\phi_t = [e^{\mu_1 t} \ e^{\mu_2 t} \ 1]^T \text{ and } \theta = [\lambda_1 \ \lambda_2 \ \eta]^T$$

and with  $e_t$  having a Gaussian distribution, the ML-estimator of  $\theta$  is then simply obtained as the linear least squares estimate. This is a direct consequence of the fact that the observations can be written in the linear regression form (4.61) which is linear in the unknown parameter  $\theta$ .

If all variables are unknown the likelihood function to be optimized becomes:

$$L = \prod_{t=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_t - \lambda_1 e^{\mu_1 t} - \lambda_2 e^{\mu_2 t} - \eta)^2}{2\sigma^2}\right). \quad (4.62)$$

Maximizing this function is equivalent to minimizing the function

$$-\log L = c + \frac{1}{2\sigma^2} \sum_{t=1}^n (y_t - \lambda_1 e^{\mu_1 t} - \lambda_2 e^{\mu_2 t} - \eta)^2. \quad (4.63)$$

Minimizing this function as a function of the unknown parameters  $\lambda_1, \lambda_2, \eta, \mu_1$  and  $\mu_2$  is an optimization problem that cannot simply be solved analytically by setting the partial derivatives equal to 0. Finding the ML estimate is now equivalent to solving a nonlinear least squares problem. More complex nonlinear optimization algorithms, as, e.g., gradient type methods, are necessary to computationally find the optimal parameter value in such cases.

## Exercises

**Exercise 4.1** Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  be i.i.d. random variables with mean  $\mu$  and variance  $\sigma^2$ .

Observe the following estimators of  $\mu$ :

$$\hat{\mu}_1 = \frac{1}{2}\mathbf{x}_1 + \frac{1}{2}\mathbf{x}_2,$$

$$\hat{\mu}_2 = \frac{2}{5}\mathbf{x}_1 + \frac{3}{5}\mathbf{x}_2,$$

$$\hat{\mu}_3 = \frac{1}{3}\mathbf{x}_1 + \frac{2}{3}\mathbf{x}_2.$$

- (a) Determine the expectation value of the three estimators.
- (b) Which of these three estimator has the smallest mean squared error (MSE)? Motivate your answer.

An estimator of the variance  $\sigma^2$  is:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

with  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$ .

- (a) Determine the expectation value of the three estimators. Is this estimator unbiased?

It can be shown that the random variables  $n\widehat{\sigma}^2/\sigma^2$  have a  $\chi^2_{n-1}$  distribution, i.e., a  $\chi^2$  with  $n - 1$  degrees of freedom.

For the expectation value and the variance of the  $\chi^2_k$  distributed variable  $y$  it holds that:  $\mathcal{E}[y] = k$  and  $var(y) = 2k$ .

- (a) Determine the variance of the estimator  $\widehat{\sigma}^2$ .

An alternative estimator of the variance  $\sigma^2$  is the so called *sample variance*:

$$\widetilde{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^2.$$

- (a) Determine which of these estimators ( $\widehat{\sigma}^2$  or  $\widetilde{\sigma}^2$ ) has the smallest mean squared error (MSE).  
 (b) Let  $\sigma/\mu = c$  and consider the following class of estimators of  $\mu$ :

$$\widehat{\mu} = \frac{a}{n} \sum_{i=1}^n \mathbf{x}_i,$$

where  $a$  is a constant. Determine the value of  $a$  (expressed in  $c$  and  $n$ ) that results in the minimum the mean squared error (MSE) of  $\widehat{\mu}$ .

**Exercise 4.2** Suppose we have a neuron that responds linearly to contrast  $x$  with slope parameter  $\theta$ , but the noise is governed by an exponential distribution:

$$f(y|x) = \alpha e^{-\alpha y}, \quad (4.64)$$

where  $\alpha = \theta x$  is the parameter of the exponential distribution (which in this case corresponds to 1 over the mean of the distribution).

Derive the maximum likelihood estimator for  $\theta$  based on a dataset of stimulus-response pairs  $\{(x_i, y_i)\}$ .

*Note that i) you should be able to find a closed form expression, just like in the Poisson and Gaussian regression problems we examined in class; ii) this is an example of a general linear model — not “generalized”, since there’s no nonlinearity between the linear stage and the noise.*

**Exercise 4.3** The binomial distribution  $Binom(n, p)$  describes the probability distribution over the number of heads from  $n$  independent coin flips, where each coin has probability  $p$  of turning up heads. We can write:

$$X \sim Binom(n, p) \quad (4.65)$$

to indicate that  $X$  is a random variable with a binomial distribution with  $n$  trials (often referred to as “ $n$  Bernoulli trials”) and parameter  $p$ . This is equivalent to saying that  $X$  has probability mass function:

$$f(X = k|n, p) = \binom{n}{k} p^k (1-p)^{n-k}. \quad (4.66)$$

- (a) Derive the maximum likelihood estimator for  $p$  given an observation of  $k$  heads from  $n$  coin flips. To do this: write out the log-likelihood,  $\log f(k|n, p)$ , take the derivative w.r.t  $p$ , set it equal to zero, and solve for  $p$ . (Show the steps in your derivation!).
- (b) You observe three binomial random variables that take on values  $k_1, k_2$ , and  $k_3$ , from  $n_1, n_2$ , and  $n_3$  coin flips. What is the maximum likelihood estimator for  $p$  given all three observations (assuming that all coins flipped had the same probability of “heads”)?

**Exercise 4.4** Let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from a distribution with probability density function

$$f(x, \theta) = \frac{1}{\theta} e^{-x/\theta}, x > 0, \theta > 0. \quad (4.67)$$

Note that, the moments of this distribution are given by

$$EX^k = \int_0^\infty \frac{x^k}{\theta} e^{-x/\theta} dx = k! \cdot \theta^k. \quad (4.68)$$

- (a) Obtain the maximum likelihood estimator of  $\theta, \hat{\theta}$ . (This should be a function of the unobserved  $x_i$  and the sample size  $n$ .) Calculate the estimate when  $x_1 = 0.50, x_2 = 1.50, x_3 = 4.00, x_4 = 3.00$ . (This should be a single number, for this dataset.)
- (b) Calculate the bias of the maximum likelihood estimator of  $\theta, \hat{\theta}$ .
- (c) Calculate the bias of the maximum likelihood estimator of maximum likelihood estimator of  $\theta, \hat{\theta}$ .

**Exercise 4.5** Let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from a distribution with probability density function

$$f(X = x|\alpha) = \alpha^{-2} x e^{-x/\alpha}, x > 0, \alpha > 0. \quad (4.69)$$

Obtain the maximum likelihood estimator of  $\alpha, \hat{\alpha}$ . Calculate the estimate when  $x_1 = 0.25, x_2 = 0.75, x_3 = 1.50, x_4 = 2.5, x_5 = 2.0$ .

**Exercise 4.6** Calculate the bias of the maximum likelihood estimator of  $\alpha, \hat{\alpha}$ .

**Exercise 4.7** Calculate the bias of the maximum likelihood estimator of maximum likelihood estimator of  $\alpha, \hat{\alpha}$ .

**Exercise 4.8** Let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from a distribution with probability density function

$$f(X = x, \beta) = \frac{1}{2\beta^3} x^2 e^{-x/\beta}, x > 0, \beta > 0. \quad (4.70)$$

Obtain the maximum likelihood estimator of  $\beta, \hat{\beta}$ . Calculate the estimate when  $x_1 = 2.00, x_2 = 4.00, x_3 = 7.50, x_4 = 3.00$ .

**Exercise 4.9** Let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from a distribution with probability density function

$$f(x, \lambda) = \lambda x^{\lambda-1}, 0 < x < 1, \lambda > 0. \quad (4.71)$$

Obtain the maximum likelihood estimator of  $\lambda$ ,  $\hat{\lambda}$ . Calculate the estimate when  $x_1 = 0.10, x_2 = 0.20, x_3 = 0.30, x_4 = 0.40$ .

**Exercise 4.10** Suppose that  $X$  is a discrete random variable with the following probability mass function: where  $0 \leq \theta \leq 1$  is a parameter:

|            |             |            |                   |                  |
|------------|-------------|------------|-------------------|------------------|
| $X = x$    | 0           | 1          | 2                 | 3                |
| $f(X = x)$ | $2\theta/3$ | $\theta/3$ | $2(1 - \theta)/3$ | $(1 - \theta)/3$ |

The following 10 independent observations were taken from such a distribution: (3,0,2,1,3,2,1,0,2,1). What is the maximum likelihood estimate of  $\theta$ ?

**Exercise 4.11** Suppose  $X_1, X_2, \dots, X_n$  are i.i.d. random variables with density function  $f(X = x|\sigma) = 1/(2\sigma) \exp(-|x|/\sigma)$ , find the maximum likelihood estimate of  $\sigma$ .

**Exercise 4.12** The Pareto distribution has been used in economics as a model for a density function with a slowly decaying tail:

$$f(X = x|x_0, \theta) = \theta x_0^\theta x^{-\theta-1}, x \geq x_0, \theta > 1. \quad (4.72)$$

Assume that  $x_0 > 0$  is given and that  $X_1, X_2, \dots, X_n$  is an i.i.d. sample. Find the MLE and CRLB of  $\theta$ .

**Exercise 4.13** Suppose that  $X_1, \dots, X_n$  form a random sample from a uniform distribution on the interval  $(0, \theta)$ , with the density function:

$$f(X = x|\theta) = \begin{cases} \frac{1}{\theta}, & \text{if } 0 \leq x \leq \theta \\ 0, & \text{otherwise} \end{cases}, \quad (4.73)$$

where the parameter  $\theta > 0$  but is unknown.

(a) Find MLE and CRLB of  $\theta$ .

(b) Now for the density function:

$$f(X = x|\theta) = \begin{cases} \frac{1}{\theta}, & \text{if } 0 < x < \theta \\ 0, & \text{otherwise} \end{cases}, \quad (4.74)$$

prove that in this case, the MLE for  $\theta$  does not exist.

**Exercise 4.14** Let  $X_1, \dots, X_n$  be an i.i.d. sample from a Poisson distribution with parameter  $\lambda$ , i.e.,

$$f(X = x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}. \quad (4.75)$$

Find the MLE and CRLB of the parameter  $\lambda$ .

**Exercise 4.15** Let  $X_1, \dots, X_n$  be an i.i.d. sample from an exponential distribution with the density function, with

$$f(X = x|\beta) = \frac{1}{\beta} e^{-\frac{x}{\beta}}, \text{ with } 0 \leq x < \infty. \quad (4.76)$$

Find the MLE and CRLB of the parameter  $\beta$ .

**Exercise 4.16** The gamma distribution has a density function as:

$$f(X = x|\alpha, \lambda) = \frac{1}{\Gamma(\alpha)} \lambda^\alpha x^{\alpha-1} e^{-\lambda x}, \text{ with } 0 \leq x < \infty. \quad (4.77)$$

Suppose that one of the parameters ( $\alpha$  or  $\lambda$ ) is known, please find the MLE of the other parameter, based on an i.i.d. sample  $X_1, \dots, X_n$ . What is the Fisher information and CRLB of these parameters?

**Exercise 4.17** Suppose that  $X_1, \dots, X_n$  form a random sample from a distribution for which the pdf  $f(X = x|\theta)$  is as follows:

$$f(X = x|\theta) = \begin{cases} \theta x^{\theta-1}, & \text{if } 0 < x \leq 1 \\ 0, & \text{otherwise} \end{cases}. \quad (4.78)$$

Also suppose that the value of  $\theta$  is unknown ( $\theta > 0$ ). Find the MLE and CRLB of  $\theta$ .





# Chapter 5

## Stochastic Processes

**After studying this chapter you can:**

- know what a stochastic process is and describe its properties via ensemble averages
- know what a Gaussian process is and know when a stochastic process is first-order and second-order stationary
- for a single (or a pair) of Wide-Sense Stationary (WSS) and zero-mean stochastic processes, you can prove the properties of its Auto- (Cross-) correlation function
- derive an estimate of the mean of ergodic WSS stochastic process that converges in the mean-square sense.
- generalize the description of a (pair of) zero-mean stochastic process(es) via its Auto- (resp. Cross-) correlation functions in the frequency domain using the Fourier Transform and prove relevant properties of the resulted spectra.

---

## 5.1 Introduction

In this chapter we start with a review of the description of random variables. This is relevant since stochastic process are (ordered) sequences of random variables. The general description of such sequences would require the definition of the Probability Density function of all individual elements in that sequence and of all their mutual combinations. This in general is extremely complicated and for the purpose of this course we make two restrictions. First we will describe stochastic processes only in terms of their so-called first order and second order statistical moments. The first order statistical moment being the mean and the second order the correlation or the covariance function. In general, unless otherwise stipulated, zero-mean stochastic processes will be considered. The second restriction is that we will assume the stochastic processes to be stationary. After a definition of this notion the important concept of Wide Sense Stationary (WSS) is introduced. In this course we will mainly be focusing on the analysis and synthesis of zero-mean WSS stochastic processes. The notion of ergodicity is introduced and its use is illustrated for the estimation of the mean value of a WSS stochastic process.

In Section 5.2.8 spectra are defined as Fourier transforms of the second order statistical moments of zero-mean WSS stochastic processes. An introduction is made on how to derive these spectra from a single and finite length realization of the stochastic process.

---

## 5.2 Stochastic Processes

Loosely speaking, stochastic or random processes are signals that depend on an independent variable, e.g. time, for which the outcome at each value of the independent variable is a random variable. In this book, we restrict the independent variable to be time and restrict to *discrete-time* random processes. These are then indexed time-sequences of a random variables. Knowing how we can characterize individual and mutual random variables, we arrive at the following definition of a stochastic process.

**Definition 5.1** (Stochastic Process). *A stochastic processes  $x(n) \in \mathbb{C}$  is an indexed time sequence of random variables,*

$$\cdots, x(-2), x(-1), x(0), x(1), x(2), \cdots$$

*with each  $x(n)$  a **random variable**. Such a sequence is characterized statistically by the individual PDF or pdf for each sample  $x(n)$  as,*

$$F_{x(n)}(\alpha) = \Pr(x(n) \leq \alpha) \quad f_{x(n)}(\alpha) = \frac{dF_{x(n)}(\alpha)}{d\alpha}$$

**and** *by the joint PDF (or pdf) for any collection of time indices  $n_1, \cdots, n_k$  as:*

$$F_{x(n_1), \dots, x(n_k)}(\alpha_1, \dots, \alpha_k) = \Pr(x(n_1) \leq \alpha_1, \dots, x(n_k) \leq \alpha_k)$$

As the result of a random experiment may provide for each random variable a particular outcome value from the sample space  $\Omega_k$ , a concatenation of the single outcomes for each random variable is called a *realization* of a stochastic process.

Let the outcome of the  $i$ th experiment with the stochastic process  $x(n)$  at time instant  $n$  be denoted by  $x_i(n)$  then one realization of this stochastic process is the discrete-time sequence:

$$\{x_i(n)\}_{n=-\infty}^{\infty}$$

---

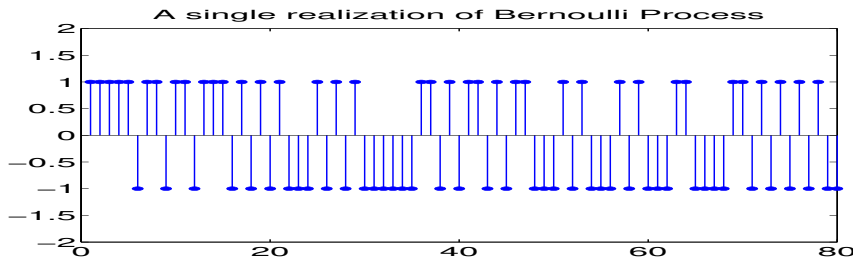
**Example 5.1 (Realization of a stochastic process)**

Let  $A$  be the outcome of throwing a “fair” die, then for a given  $\omega_0$ , we define the stochastic process  $x(n)$  as:

$$x(n) = A \cos(n\omega_0)$$

Here  $A \in \Omega = \{1, 2, \dots, 6\}$ . To each entry of  $\Omega$  we assign a probability  $\frac{1}{6}$ . Then this stochastic process is fully characterized by 6 different realizations that may occur with equal probability.

One realization of a Bernoulli process defined in Exercise 5.1 for  $p = \frac{1}{2}$  and the two possible outcomes of the Bernoulli random variable being either 1 or  $-1$ , is illustrated in Figure 5.1.



**Figure 5.1:** One realization of a Bernoulli Process with the Bernoulli random variable having the two possible outcomes 1 or  $-1$  and with a value of  $p$  (see (3.2)) equal to  $\frac{1}{2}$ .

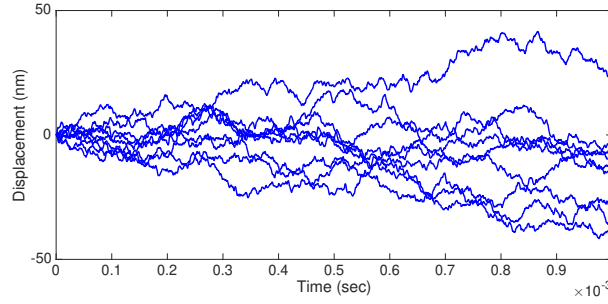
---



---

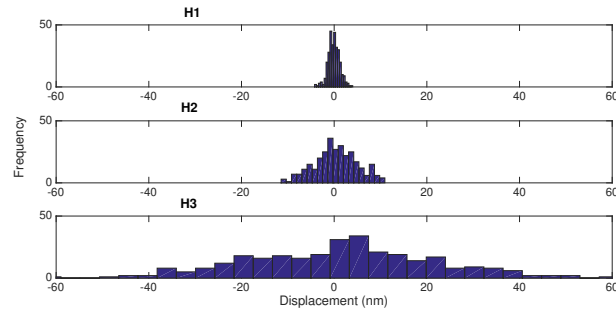
**Example 5.2 (Example 3.1 (Ct'd))**

We repeat Example 3.1 simulating 300 random experiments, each experiment observing the  $x$ -displacement of a particle in a water solution for a period of  $1\mu s$ . Each random experiment starts from the same initial conditions. From these 300 experiments we plot in Figure 5.2 ten different realizations. We clearly observe the randomness or non-repeatability of the experiment whereby each new experiment or trial delivers a different curve of the  $x$  displacement. Subsequently, these 300 realizations determine the possible outcomes of the random variables that result by considering the  $x$ -position of a particle at one particular time instant. For example considering the time instances  $(0.005\mu s, 0.05\mu s, 1\mu s)$ , we used the resulting outcomes to plot in Figure 5.3 their corresponding Histograms.



**Figure 5.2:** Ten realization of  $x$ -position of a particle as simulated in Example 3.1.

When considering the histogram as an approximation of the pdf of the random variable, we may conclude from this figure that the statistical properties of the 3 considered random variables is different.



**Figure 5.3:** The Histograms of the Random Variable taken as the  $x$ -position of a particle as simulated in Example 3.1 at different time instances ( $0.005\mu s, 0.05\mu s, 1\mu s$ ) for 300 realizations.

### 5.2.1 Ensemble Averages for a Stochastic Process

For each time index  $n$  the quantity  $x(n)$  of a stochastic process is a random variable. Therefore we can generalize the notions of Ensemble averages given in Table 3.2 now as in Table 5.1. The Ensemble notions turn out to be functions. But generally we may drop the part function in referring to the particular Ensemble Average.

| Quantity                     | Symbol          | Definition                                  |
|------------------------------|-----------------|---|
| Mean (Function):             | $m_x(n)$        | $E[x(n)]$                                   |
| Variance (Function):         | $\sigma_x^2(n)$ | $E[ x(n) - m_x(n) ^2]$                      |
| Auto-correlation (Function): | $r_x(k, \ell)$  | $E[x(k)x(\ell)^*]$                          |
| Auto-covariance (Function):  | $c_x(k, \ell)$  | $E[(x(k) - m_x(k))(x(\ell) - m_x(\ell))^*]$ |

**Table 5.1:** Different Moments for a stochastic process  $x(n)$

## 5.2.2 Ensemble Averages for two Stochastic Processes

The different Ensemble Averages for a single Stochastic Process can easily be generalized for two stochastic processes. Let these two processes be denoted by  $x(n)$  and  $y(n)$  resp. then we summarize different Ensemble averages for two Stochastic Processes in Table 5.2.

| Quantity                      | Symbol            | Definition                                  |
|-------------------------------|-------------------|---|
| Mean (Function):              | $m_x(n)$          | $E[x(n)]$                                   |
|                               | $m_y(n)$          | $E[y(n)]$                                   |
| Cross-correlation (Function): | $r_{xy}(k, \ell)$ | $E[x(k)y(\ell)^*]$                          |
| Cross-covariance (Function):  | $c_{xy}(k, \ell)$ | $E[(x(k) - m_x(k))(y(\ell) - m_y(\ell))^*]$ |

**Table 5.2:** Different Moments for two stochastic processes  $x(n)$  and  $y(n)$ .

Based on the definition of the ensemble averages between two stochastic process, we can now generalize the notions of uncorrelatedness and orthogonality defined in Section 3.2.4.

Two stochastic processes  $x(n)$  and  $y(n)$  are *uncorrelated* when their Cross-correlation function satisfies:

$$r_{xy}(k, \ell) = m_x(k)m_y(\ell)^* \quad (5.1)$$

Based on the relationship between the Cross-correlation function  $r_{xy}(k, \ell)$  and the Cross-covariance function given as:

$$c_{xy}(k, \ell) = r_{xy}(k, \ell) - m_x(k)m_y(\ell)^*$$

we can conclude that for uncorrelated stochastic processes the cross-covariance function is zero.

Two stochastic processes  $x(n)$  and  $y(n)$  are *orthogonal* when their Cross-correlation function satisfies:

$$r_{xy}(k, \ell) = 0 \quad (5.2)$$

Therefore uncorrelated stochastic processes are orthogonal provided the mean of one of them is equal to zero.

In these course notes, we restrict the description of stochastic processes to the Ensemble averages defined in Tables 5.1 and 5.2. Further we assume throughout the book, unless otherwise indicated, that all stochastic processes have **mean zero**.

## 5.2.3 Gaussian Processes

The definition of two jointly Gaussian random variables may be generalized for a stochastic process. For that purpose, we collect  $n$  random variables of a real, stochastic process  $x(n)$  into the vector  $\mathbf{x}$ ,

$$\mathbf{x} = [x(1) \ x(2) \ \cdots \ x(n)]^T$$

then  $\mathbf{x}$  is a Gaussian Random vector with the non-singular, covariance matrix given as:

$$C_x = E \left[ \begin{bmatrix} x(1) - E[x(1)] \\ \vdots \\ x(n) - E[x(n)] \end{bmatrix} \begin{bmatrix} (x(1) - E[x(1)]) & \cdots & (x(n) - E[x(n)]) \end{bmatrix} \right]$$

and mean vector given as:

$$\mathbf{m}_x^T = [E[x(1)] \quad \cdots \quad E[x(n)]]$$

when its joint probability density function has the following form:

$$f_{\mathbf{x}}(\boldsymbol{\alpha}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\det(C_x)|^{\frac{1}{2}}} \exp^{-\frac{1}{2}(\boldsymbol{\alpha} - \mathbf{m}_x)^T C_x^{-1} (\boldsymbol{\alpha} - \mathbf{m}_x)} \quad (5.3)$$

A stochastic process  $x(n)$  is now *Gaussian* if every **finite** collection of samples of  $x(n)$  are jointly Gaussian. As we conclude from (5.3) the pdf of such a finite collection of samples is fully determined by its mean and its covariance matrix.

Gaussian Processes are an important tool to analyse and describe complex processes. They are widely used in signal analysis [4]. This is because in many cases real-life data can very well be approximated by Gaussian processes. This is a consequence of the Central Limit Theorem [5].

## 5.2.4 Stationary Processes

The characterization of a stochastic process as in Definition 5.1 via the individual and joint PDFs (or pdfs) is quite complicated. Even when restricting to the description via first and second moment Ensemble averages as outlined in Section 5.2.1 requires the definition of time varying functions of one or two variables. In practice however, these characteristic statistical quantities are not dependent explicitly on time. The notion of time-invariance in the context of stochastic process is called *stationarity*. Such notion will turn out to be of great help when estimating statistical characteristics, like the mean of stochastic process, from a single realization.

Different types of stationarity have been used in the literature.

### First-order Stationarity

**Definition 5.2** (First-order Stationarity). *A stochastic process  $x(n)$  with individual pdf  $f_{x(n)}(\alpha)$  as defined in Definition 5.1 is first-order stationary if and only if*

$$f_{x(n)}(\alpha) = f_{x(n+k)}(\alpha) \quad \forall k$$

This definition leads to the following corollary.

**Corollary 5.3.** *Let a stochastic process  $x(n)$  with individual pdf  $f_{x(n)}(\alpha)$  as defined in Definition 5.1 be first-order stationary, then its mean is constant.*

*Proof.* Exercise 5.7 calls for a proof of this Corollary. □

## Second-order Stationarity

**Definition 5.4** (Second-order Stationarity). *A stochastic process  $x(n)$  with mutual pdf  $f_{x(n_1),x(n_2)}(\alpha_1, \alpha_2)$  is second-order stationary if and only if*

$$f_{x(n_1),x(n_2)}(\alpha_1, \alpha_2) = f_{x(n_1+k),x(n_2+k)}(\alpha_1, \alpha_2) \quad \forall k$$

This definition leads to the following corollary.

**Corollary 5.5.** *Let a stochastic process  $x(n)$  with mutual pdfs  $f_{x(n_1),x(n_2)}(\alpha_1, \alpha_2)$  be second-order stationary, then its Auto-correlation function  $r_x(k, \ell)$  only depends on the difference  $k - \ell$ , i.e.*

$$r_x(k, \ell) = r_x(k - \ell, 0) := r_x(k - \ell) \quad \forall k, \ell$$

The difference  $k - \ell$  in the above Auto-correlation function is referred to as the *lag* of the Auto-correlation function.

*Proof.* Exercise 5.7 calls for a proof of this Corollary. □

### 5.2.5 Wide Sense Stationary (WSS) Processes

The main notion of stationarity of practical interest that will be mainly used throughout this book is the notion of *Wide Sense Stationarity*, abbreviated as WSS. This is a weaker form of stationarity that only imposes conditions on the ensemble averages (first and second moment).

#### For a single Stochastic Process

The notion of WSS for a single stochastic process is defined in the following definition.

**Definition 5.6** (Wide-Sense Stationarity). *A stochastic process  $x(n)$  is Wide-Sense Stationary (WSS) if it satisfies the following three criteria,*

1.  $m_x(k) = m_x < \infty$
2.  $r_x(k, \ell) = r_x(k - \ell) \quad \forall k, \ell$
3.  $c_x(0) < \infty$ , i.e. its variance is finite.

**Remark 5.7** (Constant variance). *Here variance is also not time-dependent (constant). It can be proven in the following way, using the properties that the mean of a WSS process is constant and auto-correlation function depends only on the difference  $k - l$ :*

$$\begin{aligned} \text{var} &= E[(x(n) - E[x(n)])^2] = E[(x(n) - m_x)^2] \\ &= E[x(n)^2 - 2x(n)m_x + m_x^2] = E[x(n)x(n)] - m_x^2 \\ &= r_x(0, 0) - m_x^2 = r_x(0) - m_x^2 \end{aligned}$$

$r_x(0)$  is not time-dependent,  $m_x$  is a constant.

The fact that it is a weaker notion stems for example from the fact that the first condition follows from first-order stationarity. However the reverse does not need to be true.

Stochastic processes that are WSS have a number of interesting properties that make them useful from a practical perspective.

**Property 5.8 (Symmetry).** *If  $x(n)$  is a complex stochastic process that is WSS, then its Auto-correlation function is conjugate symmetric, i.e.*

$$r_x(k) = r_x^*(-k)$$

*If  $x(n)$  is real and WSS, then its Auto-correlation function is symmetric, i.e.*

$$r_x(k) = r_x(-k)$$

**Property 5.9 (Non-negative Variance).** *If  $x(n)$  is a complex stochastic process that is WSS, then its Variance is non-negative, i.e.*

$$r_x(0) = E[|x(n)|^2] \geq 0$$

**Property 5.10 (Maximum).** *If  $x(n)$  is a complex stochastic process that is WSS, then its Auto-correlation function reaches its maximum value for lag 0, i.e.*

$$r_x(0) \geq |r_x(k)| \quad \forall k$$

## For two Stochastic Processes

The notion of WSS for two stochastic processes is defined in the following definition.

**Definition 5.11 (Wide-Sense Stationarity).** *Two stochastic processes  $x(n)$  and  $y(n)$  are jointly Wide-Sense Stationary (WSS) if it satisfies the following criteria,*

1.  $x(n)$  and  $y(n)$  are WSS
2.  $r_{xy}(k, \ell) = r_{xy}(k - \ell) \quad \forall k, \ell$

The Cross-correlation function possesses the following three properties.

**Property 5.12.** *If  $x(n)$  and  $y(n)$  are complex stochastic processes and they are jointly WSS, then its Cross-correlation function  $r_{xy}(k)$  satisfies,*

$$r_{xy}(k) = r_{yx}^*(-k)$$

*(be aware of the change of order in the variables in the subscript definition of the Cross-correlation function). If both  $x(n)$  and  $y(n)$  are real, then,*

$$r_{xy}(k) = r_{yx}(-k)$$

**Property 5.13.** *If  $x(n)$  and  $y(n)$  are complex stochastic processes and they are jointly WSS, with Auto-correlation function resp. equal to  $r_x(k)$  and  $r_y(k)$ , then its Cross-correlation function  $r_{xy}(k)$  satisfies,*

$$|r_{xy}(k)| \leq \sqrt{r_x(0)r_y(0)}$$



**Property 5.14.** If  $x(n)$  and  $y(n)$  are complex stochastic processes and they are jointly WSS, with Auto-correlation function resp. equal to  $r_x(k)$  and  $r_y(k)$ , then its Cross-correlation function  $r_{xy}(k)$  satisfies,

$$|\operatorname{Re}(r_{xy}(k))| \leq \frac{r_x(0) + r_y(0)}{2}$$

### 5.2.6 Autocorrelation matrix of a Stochastic Process

When  $x(n)$  is a complex WSS stochastic process with Auto-correlation function  $r_x(k)$ , an important notion that will be used later in Chapter 9 is the Auto-correlation matrix. We restrict in this chapter to the case of a  $3 \times 3$  matrix. For this case, the Auto-correlation matrix is defined as the covariance matrix of the

vector  $\mathbf{x} = \begin{bmatrix} x(n) \\ x(n+1) \\ x(n+2) \end{bmatrix}$  and is defined as:

$$R_{\mathbf{x}} = E[\mathbf{x}\mathbf{x}^H] = E \left[ \begin{bmatrix} x(n) \\ x(n+1) \\ x(n+2) \end{bmatrix} \begin{bmatrix} x^*(n) & x^*(n+1) & x^*(n+2) \end{bmatrix} \right] \quad (5.4)$$

The Auto-correlation matrix can be expressed in terms of the Auto-correlation function as follows:

$$R_x = \begin{bmatrix} r_x(0) & r_x^*(1) & r_x^*(2) \\ r_x(1) & r_x(0) & r_x^*(1) \\ r_x(2) & r_x(1) & r_x(0) \end{bmatrix} \quad (5.5)$$

### 5.2.7 Ergodicity

The description of random variables and stochastic processes via Ensemble Averages is practically relevant. For that reason it is of interest to try to estimate Ensemble Averages from experiments. For random variables one can estimate Ensemble Averages from the experimentally retrieved outcomes. For example let  $x_i$  denote the experimental outcome of the number of eyes when throwing a “fair” dye, and when we repeat this experiment  $L$  times, then the mean value of the random variable “number of eyes” could be provided by the average of the form,

$$\hat{m}_x = \frac{1}{L} \sum_{i=1}^L x_i$$

When we have access to multiple realizations generated by multiple random experiments with a stochastic process, we could repeat the above averaging process. Consider for example the 300 realizations in Example 5.2, we could aim at estimating the mean at time instance  $0 \leq t \leq 1 \mu s$  via a similar average,

$$\hat{m}_x(t) = \frac{1}{300} \sum_{i=1}^{300} x_i(t) \quad (5.6)$$

However when the experimental conditions hamper the retrieval of multiple realizations of the same stochastic process, one is interested in getting estimates of Ensemble Averages from a single realization. It is here that the notion of ergodicity comes in.

*Ergodicity* is loosely speaking the condition that allows to retrieve Ensemble Averages of a stochastic process from a single realization. In this book we restrict ourselves to the mean and Auto-correlation function Ensemble Averages.

### Ergodicity in the mean

Since use will be made of a single realization of a stochastic process  $x(n)$  we denote that realization also as the time series  $x(n)$ . When  $N$  samples of a realization are available it is denoted as:  $\{x(n)\}_{n=0}^{N-1}$ . Using these samples we could propose the following *time average*:

$$\hat{m}_x(N) = \frac{1}{N} \sum_{n=0}^{N-1} x(n) \quad (5.7)$$

Be aware of the fact that we now sum over the time index  $n$  (and not the realization index  $i$  as in the example given in (5.6) . We are now interested in which properties of the stochastic process, the time average (5.7) becomes equal to the Ensemble mean of the stochastic process. As (5.7) is a single number, it means that we have at least to consider the stochastic process having a constant mean.

To rephrase our interest more precisely, we are interested in the property summarized in the following Definition.

**Definition 5.15** (Ergodic in the mean). *When  $x(n)$  is a WSS stochastic process, then  $x(n)$  is ergodic in the mean, if the time average  $\hat{m}_x(N)$  in (5.7) converges in the mean square sense [5], given as,*

$$\lim_{N \rightarrow \infty} E[|\hat{m}_x(N) - m_x|^2] = 0 \quad (5.8)$$

The convergence in the mean square sense is also written as,

$$\lim_{N \rightarrow \infty} \hat{m}_x(N) = m_x \quad \text{or} \quad \Pr \left[ \lim_{N \rightarrow \infty} \hat{m}_x(N) = m_x \right] = 1$$

A condition on the Auto-correlation function is provided in the following Theorem.

**Theorem 5.16** (Condition to guarantee Ergodicity in the mean). *A WSS stochastic process  $x(n)$  is ergodic in the mean if and only if its Auto-covariance function  $c_x(k)$  satisfies,*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=-N+1}^{N-1} \left(1 - \frac{|k|}{N}\right) c_x(k) = 0$$

*Proof.* Following Definition 5.15 for ergodicity in the mean we consider the following variance,

$$\text{Var}(\hat{m}_x(N)) = E[|\hat{m}_x(N) - m_x|^2]$$

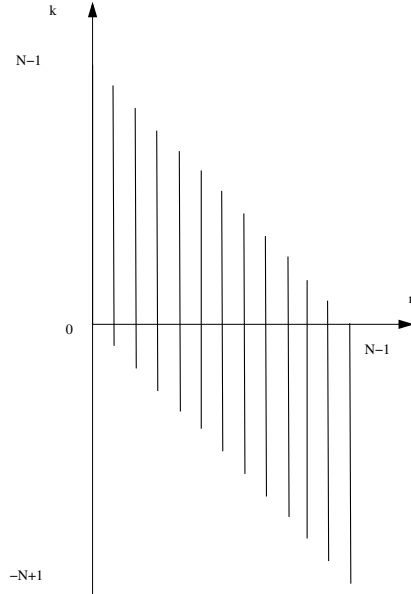
Using the definition in (5.7), and the definition of the Auto-covariance function  $c_x(k)$ , this variance  $\text{Var}(\hat{m}_x(N))$  can be written as,

$$\begin{aligned}\text{Var}(\hat{m}_x(N)) &= E\left[\left|\frac{1}{N} \sum_{n=0}^{N-1} (x(n) - m_x)\right|^2\right] \\ &= \frac{1}{N^2} \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} E[(x(m) - m_x)(x(n) - m_x)^*] \\ &= \frac{1}{N^2} \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} c_x(m - n)\end{aligned}$$

With a change of variables defining  $k = m - n$ , we have,

$$= \frac{1}{N^2} \sum_{n=0}^{N-1} \sum_{k=-n}^{N-1-n} c_x(k) \quad (5.9)$$

Now we consider the indices  $n$  and  $k$  as displayed in Figure 5.4. From (5.9) we



**Figure 5.4:** The two indices  $k$  and  $n$  in double sum (5.9).

see that for  $n = 0$ , the variable  $k$  runs from 0 to  $N - 1$ , and for  $n = 1$   $k$  runs from  $-1$  to  $N - 2$ . The line segment drawn at  $n = 1$  runs from  $k = -1$  to  $N - 2$ . This is used to indicate the set of integers  $k$  on that line segment.

Similarly we can interpret the other line segments in Figure 5.4. If we now take the intersections of this line segments for  $k = 0$ , we see that we have  $c_x(0)$   $N$  times. For  $k = 1$  we have  $c_x(1)$   $N - 1$  times and for  $k = -1$  we have  $c_x(-1) =$

$c_x(1)$  also  $N - | - 1 |$  times. Continuing in this way, (5.9) equals,

$$\begin{aligned} &= \frac{1}{N^2} \sum_{k=-N+1}^{N-1} (N - |k|) c_x(k) \\ &= \frac{1}{N} \sum_{k=-N+1}^{N-1} \left(1 - \frac{|k|}{N}\right) c_x(k) \end{aligned}$$

and the proof is completed.  $\square$

In [5] a more refined condition on the Auto-covariance function is given. There the conditions for a WSS stochastic process  $x(n)$  with Auto-covariance function  $c_x(k)$  to be ergodic in the mean are that,

1.  $c_x(0) < \infty$  and
2.  $\lim_{k \rightarrow \infty} c_x(k) = 0$ .

This shows that for ergodicity in the mean the variance of the stochastic process need to be finite and the Auto-covariance function should go to zero for large lags.

### Auto-correlation Ergodic

When a single realization of a WSS stochastic process  $x(n)$  is available in the form of the samples,

$$\{x(n)\}_{n=-k}^{N-1}$$

So the length of the sequence depends on the index  $k$  (and  $N$ ). Based on this realization we can define a time average of the Auto-correlation function as,

$$\hat{r}_x(k, N) = \frac{1}{N} \sum_{n=0}^{N-1} x(n) x^*(n - k) \quad (5.10)$$

When this time average converges to the true Auto-correlation function  $r_x(k)$  the WSS stochastic process  $x(n)$  is called *Auto-correlation Ergodic*.

This condition stipulates a condition on the Auto-covariance function as summarized in the following Theorem [2], which we state without proof.

**Theorem 5.17** (Auto-correlation Ergodic). [2] A WSS Gaussian process  $x(n)$  is Auto-correlation ergodic if and only if its Auto-covariance function  $c_x(k)$  satisfies,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} c_x^2(k) = 0$$

### 5.2.8 WSS Processes in the Frequency Domain

As we assume WSS to have mean zero, the description of a WSS process in the frequency domain calls for considering the Fourier transforms of the Auto-correlation function for a single stochastic process and of the Cross-correlation function for two stochastic processes.

## Power Spectrum

The Fourier transform of the Auto-correlation is considered in the following Definition.

**Definition 5.18** (Power Spectrum). *The Power Spectrum (or power spectral density) of a WSS stochastic process  $x(n)$  is the DTFT of its Auto-correlation function  $r_x(k)$ :*

$$P_x(e^{j\omega}) = \sum_{k=-\infty}^{\infty} r_x(k) e^{-j\omega k}$$

Using the definition of the inverse Fourier transform as given in (2.10), we can retrieve the Auto-correlation function back from the Power Spectrum as,

$$r_x(k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} P_x(e^{j\omega}) e^{j\omega k} d\omega \quad (5.11)$$

As done in Chapter 2, we also make use of the z-transform of the Auto-correlation function  $r_x(k)$ .

$$P_x(z) = \sum_{k=-\infty}^{\infty} r_x(k) z^{-k} \quad (5.12)$$

We will both refer to  $P_x(e^{j\omega})$  and  $P_x(z)$  as the Power Spectrum of  $x(n)$ , as the use of either the argument  $e^{j\omega}$  (with  $\omega \in \mathbb{R}$  indicating the frequency dependency) or  $z \in \mathbb{C}$  will be clear from the context. The use of the argument  $z$  is often preferred when doing mathematical calculations. The result can easily be transformed to the argument  $e^{j\omega}$  by taking  $z = e^{j\omega}$ .

The Power spectrum has the following 4 properties.

**Property 5.19** (real and Symmetry). *If  $x(n) \in \mathbb{C}$  is a complex stochastic process that is WSS, then its Power Spectrum satisfies,*

$$P_x(e^{j\omega}) = P_x^*(e^{j\omega})$$

*that is  $P_x(e^{j\omega}) \in \mathbb{R}$  and*

$$P_x(z) = P_x^*(1/z^*)$$

*If  $x(n) \in \mathbb{R}$ , then,*

$$P_x(e^{j\omega}) = P_x(e^{-j\omega})$$

*that is  $P_x(e^{j\omega})$  is even and*

$$P_x(z) = P_x^*(z^*)$$

**Property 5.20** (Positivity of Spectrum). *If  $x(n) \in \mathbb{C}$  is a complex stochastic process that is WSS, then its Power Spectrum satisfies,*

$$P_x(e^{j\omega}) \geq 0$$

**Property 5.21** (Total Power). *If  $x(n) \in \mathbb{C}$  is a complex stochastic process that is WSS, then its Power Spectrum satisfies,*

$$E[|x(n)|^2] = \frac{1}{2\pi} \int_{-\pi}^{\pi} P_x(e^{j\omega}) d\omega$$

**Remark 5.22** (Power in the Signal). *The integral under the Power spectrum represents the Total power of that signal. When the integration is only done over a subinterval it represents the Power of the signal in that interval only.*

**Property 5.23** (Periodic Component). *If  $x(n) \in \mathbb{C}$  is the sum of a complex stochastic process that is WSS and an harmonic signal, then its Power Spectrum contains unit pulse(s).*

Exercise 5.5 calls for the calculation of the Power Spectrum of a stochastic process containing Harmonics.

Related to Property 5.23 a stochastic process that does not contain harmonics, is called a *regular* stochastic process.

### Asymptotic Estimate of the Power Spectrum

First we are trying to propose an estimate for the Power Spectrum. For that purpose we inspect its definition more closely and study the effect of having only a finite number of samples of a single realization of the stochastic process  $x(n)$ . Namely we will consider the following finite series to be given:

$$\{x(n)\}_{n=-N}^N \quad (5.13)$$

For the sequence  $x(n)$  we define the following Fourier transforms, assuming that it exists,

$$X(e^{j\omega}) = \sum_{k=-\infty}^{\infty} x(k)e^{-j\omega k}$$

and if we consider only a finite number of samples as indicated in (5.13), the following approximation of the Fourier transform is defined,

$$X_N(e^{j\omega}) = \sum_{k=-N}^N x(k)e^{-j\omega k} \quad (5.14)$$

Using these Fourier transforms we are now in a position, to derive an approximate of the Power spectrum using a single finite length realization. This approximation will be called the *Periodogram*.

To derive the approximation recall the definition of the Power Spectrum in Definition 5.18 and that of the Auto-correlation function as,

$$P_x(e^{j\omega}) = \sum_{k=-\infty}^{\infty} r_x(k)e^{-j\omega k} = \sum_{k=-\infty}^{\infty} E[x(n)x^*(n-k)]e^{-j\omega k}$$

Since  $E[\cdot]$  is a linear operator and using the Time Shift property of the DTFT in Table 2.2, the above is equal to,

$$\begin{aligned} P_x(e^{j\omega}) &= E\left[x(n) \sum_{k=-\infty}^{\infty} x^*(n+k)e^{j\omega k}\right] \\ &= E\left[x(n) \left[ \sum_{k=-\infty}^{\infty} x(n+k)e^{-j\omega k} \right]^*\right] \\ &= E\left[x(n)e^{-j\omega n} X^*(e^{j\omega})\right] \end{aligned}$$

We can now consider  $x(n)e^{-j\omega n}X^*(e^{j\omega})$  as a complex stochastic process  $x(n)e^{-j\omega n}$  scaled by the complex number  $X^*(e^{j\omega})$ . Therefore the mean value  $E[x(n)e^{-j\omega n}X^*(e^{j\omega})]$  can be approximated as,

$$E[x(n)e^{-j\omega n}X^*(e^{j\omega})] \approx \frac{1}{2N+1} \sum_{n=-N}^N x(n)e^{-j\omega n}X^*(e^{j\omega})$$

Substituting in this equation the approximation of the Fourier transform  $X(e^{j\omega})$  given in (5.14), we obtain the following approximation of the Power Spectrum,

$$P_x(e^{j\omega}) \approx \frac{1}{2N+1} X_N(e^{j\omega})X_N^*(e^{j\omega}) = \frac{1}{2N+1} |X_N(e^{j\omega})|^2$$

This approximation is called the *Periodogram*. It is referred to as  $\hat{P}_x(e^{j\omega})$ ,

$$\hat{P}_x(e^{j\omega}, N) = \frac{1}{2N+1} |X_N(e^{j\omega})|^2 \quad (5.15)$$

The periodogram is a widely used estimate of the Power spectrum as it makes use of the Fourier transform. Especially when using the so-called fast Fourier transform [2], it allows to handle large data sets in an efficient manner.

The ergodicity result for the periodogram is provided in the following Theorem, stated without proof.

**Theorem 5.24** (Asymptotic Estimate of the Power Spectrum). *Let the Auto-correlation function of a zero-mean, WSS Gaussian process  $x(n)$  satisfy,*

$$\sum_{k=-\infty}^{\infty} |k|r_x(k) < \infty$$

*and let the Periodogram be defined as in (5.15), then the mean of the random variable  $\hat{P}_x(e^{j\omega}, N)$  asymptotically as the number of data points  $N$  go to  $\infty$  equals to the true Power Spectrum. This is indicated as,*

$$P_x(e^{j\omega}) = \lim_{N \rightarrow \infty} E[\hat{P}_x(e^{j\omega}, N)]$$

It is important to note that the expectation  $E[\hat{P}_x(e^{j\omega}, N)]$  is taken for a fixed  $\omega$ . Hence in order to provide a “good” estimate of the Power spectrum, not only  $\hat{P}_x(e^{j\omega}, N)$  for large  $N$  is necessary, but also a good estimate of the Ensemble average  $E[\hat{P}_x(e^{j\omega}, N)]$ . This is why good estimator of the spectrum also try to derive multiple Periodograms from a single finite time length realization. For a discussion on such estimator we refer to Chapter 8 of [2].

## Cross Spectrum

The Fourier transform of the Cross-correlation between two WSS stochastic processes is considered in the following Definition.

**Definition 5.25** (Cross Spectrum). *The Cross Spectrum (or cross spectral density) between two jointly WSS stochastic processes  $x(n), y(n)$  is the DTFT of their Cross-correlation function  $r_{xy}(k)$ :*

$$P_{xy}(e^{j\omega}) = \sum_{k=-\infty}^{\infty} r_{xy}(k)e^{-j\omega k}$$

Using the definition of the inverse Fourier transform as given in (2.10), we can retrieve the Cross-correlation function back from the Cross Spectrum as,

$$r_{xy}(k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} P_{xy}(e^{j\omega})e^{j\omega k} d\omega \quad (5.16)$$

As done in Chapter 2, we also make use of the z-transform of the Cross-correlation function  $r_{xy}(k)$ .

$$P_{xy}(z) = \sum_{k=-\infty}^{\infty} r_{xy}(k)z^{-k} \quad (5.17)$$

We will both refer to  $P_{xy}(e^{j\omega})$  and  $P_{xy}(z)$  as the Cross Spectrum between  $x(n)$  and  $y(n)$ . Remark that the order of the indices  $xy$  is relevant in the definition 5.25 and in (5.17).

The Cross Spectrum has the following 4 properties.

**Property 5.26** (Complex). *If  $x(n), y(n) \in \mathbb{R}$  are jointly WSS, then its Cross Spectrum is complex,*

$$P_{xy}(e^{j\omega}) \in \mathbb{C}$$

**Property 5.27** (Even and Odd Parts of Cross Spectrum). *If  $x(n), y(n) \in \mathbb{R}$  are jointly WSS, then,*

$$P_{xy}(e^{j\omega}) = P_{xy}^*(e^{-j\omega})$$

*That is the Real part of the complex number  $P_{xy}(e^{j\omega})$  is an even function, while the Imaginary part is an odd function.*

**Property 5.28.** *If  $x(n), y(n) \in \mathbb{C}$  are jointly WSS, then its Cross Spectrum satisfies,*

$$P_{xy}(z) = P_{yx}^*(1/z^*)$$

*If instead  $x(n), y(n) \in \mathbb{R}$ , then,*

$$P_{xy}(e^{j\omega}) = P_{yx}(e^{-j\omega})$$

**Property 5.29** (Orthogonal). *If  $x(n), y(n) \in \mathbb{R}$  are jointly WSS, and orthogonal then its Cross Spectrum satisfies,*

$$P_{xy}(e^{j\omega}) = 0$$



---

## References

- [1] A. Leon-Garicia, *Probability and Random Processes for Electrical Engineering*. Addison-Wesley, 1994 (2nd-edition).
  - [2] M. H. Hayes, *Statistical Digital Signal Processing and Modeling*. New York: John Wiley and Sons, 1996.
  - [3] A. Einstein, *A New Determination of Molecular Dimensions*, PhD Thesis, University of Zurich, 1905.
  - [4] C.E. Rasmussen and C.K. I. Williams, *Gaussian Processes for Machine Learning*, MIT Press, 2006.
  - [5] A. Papoulis and S. Unnikrishna Pillai, *Probability, Random Variables, and Stochastic Processes*, Mc Graw Hill, 4rd Edition, 2002.
- 

## Exercises

**Exercise 5.1** Let  $a_k$  be a Bernoulli random variable for each value of  $k$ , with outcomes  $+1$  and  $-1$  as defined in Definition 3.2 for  $p$  equal to  $\frac{1}{2}$ , and let  $E[a_k a_j] = 0 \forall k, j : k \neq j$ , then we define the Bernoulli process  $x(n)$  as follows:

$$x(n) = \sum_{k=-\infty}^{\infty} a_k \delta(n - k)$$

Determine,

- (a)  $E[x(n)]$ , and,
- (b)  $E[x(n_1)x(n_2)]$ .

**Exercise 5.2** Let  $x, y$  be two random variables with a Gaussian distribution, then show that when the correlation coefficient  $\rho_{xy}$  is zero, that the two random variables are *independent*.

**Exercise 5.3** If  $x$  is Gaussian random variable, with moments defined in Table 3.2, then prove that

$$E[x^n] = \begin{cases} 1 \times 3 \times 5 \times \cdots \times (n-1) \sigma_x^n & : n \text{ even} \\ 0 & : n \text{ odd} \end{cases}$$

**Exercise 5.4** If  $x$  and  $y$  are two random variables with a Gaussian distribution, then the optimum nonlinear estimator for  $y$  from  $x$  given as:

$$\hat{y} = g(x)$$

that minimizes the mean-square error  $E[(y - g(x))^2]$  is the linear estimator given as:

$$\hat{y} = ax + b$$

Determine the values of  $a$  and  $b$  in terms of the moments given in Table 3.2.

**Exercise 5.5** Let  $\phi$  be a Random Variable with pdf,

$$f_{\phi}(\alpha) = \begin{cases} \frac{1}{2\pi} & -\pi \leq \alpha < \pi \\ 0 & \text{otherwise} \end{cases}$$

further let  $A$  and  $\omega_0$  be fixed constants ( $\in \mathbb{R}$ ), then we define the stochastic process  $x(n)$  as:

$$x(n) = A \sin(n\omega_0 + \phi)$$

Determine,

- (a)  $m_x(n)$ , and,
- (b)  $r_x(k, \ell)$ .
- (c) The Power Spectrum  $P_x(e^{j\omega})$  as the Fourier transform of its Auto-correlation function (if it exists).

**Exercise 5.6** Let the stochastic process  $z(n)$  be the sum of two uncorrelated stochastic processes  $x(n)$  and  $y(n)$ , then show that:

$$r_z(k, \ell) = r_x(k, \ell) + r_y(k, \ell)$$

**Exercise 5.7** Prove Corollaries 5.3 and 5.5.

**Exercise 5.8** Prove Properties 5.8, 5.9 and 5.10 of a WSS stochastic process  $x(n)$ .

**Exercise 5.9** Prove Properties 5.12, 5.13 and 5.14 of two jointly WSS stochastic processes  $x(n)$  and  $y(n)$ .

**Exercise 5.10** Let  $v_1(n)$  and  $v_2(n)$  be zero-mean, uncorrelated jointly WSS stochastic processes, and let for constant values of  $\omega$  and  $\phi$  the following stochastic processes be defined:

$$x(n) = \sin(\omega n) + v_1(n) \quad y(n) = \sin \omega n + \phi + v_2(n)$$

then,

- (a) Determine the Cross-correlation function  $r_{xy}(k)$ .
- (b) Is the pair of stochastic processes jointly WSS? Motivate your answer.
- (c) Determine a procedure to retrieve the phase difference  $\phi$  from the Cross-correlation function  $r_{xy}(k)$ .

**Exercise 5.11** For a complex WSS stochastic process with Auto-correlation function  $r_x(k)$ , show that the Auto-covariance matrix  $R_x$  defined in Section 5.2.6 is equal to,

$$R_{\mathbf{x}} = E \left[ \begin{bmatrix} x^*(n) \\ x^*(n-1) \\ x^*(n-2) \end{bmatrix} \begin{bmatrix} x(n) & x(n-1) & x(n-2) \end{bmatrix} \right]$$

**Exercise 5.12** For a complex WSS stochastic process with Auto-correlation function  $r_x(k)$ , show that the Auto-covariance matrix  $R_x$  defined in Section 5.2.6, show that this matrix satisfies,

- (a)  $R_x$  is Hermitian Toeplitz.
- (b)  $R_x \geq 0$  and therefore the eigenvalues of  $R_x$  are real and non-negative.

**Exercise 5.13** Prove the Property 5.19.

**Exercise 5.14** Prove the Property 5.20.

**Exercise 5.15** Prove the Property 5.21.

**Exercise 5.16** Prove the Property 5.26.

**Exercise 5.17** Prove the Property 5.27.

**Exercise 5.18** Prove the Property 5.28.

**Exercise 5.19** Prove the Property 5.29

**Exercise 5.20** Let  $u$  be a uniform random variable on  $[0, 1]$ .

- (a) Compute the expected value of  $u$ .
- (b) Compute the variance of  $u$ .
- (c) Let  $x$  and  $y$  be random variables defined as

$$\begin{aligned}x &= \cos(2\pi u) \\y &= \sin(2\pi u)\end{aligned}$$

Are  $x$  and  $y$  orthogonal?

- (d) Are  $x$  and  $y$  uncorrelated?
- (e) Do you expect  $x$  and  $y$  to be independent? Please motivate your answer.

**Exercise 5.21** A Poisson process is a counting process where events are counted over time. Events can be anything ranging from earthquakes, failures in a power plant, shooting electrons, etc. For a Poisson process  $x(t)$  with expected value  $\lambda t$ , the probability of  $k$  events in a time interval of length  $t$  is given by

$$Pr(x(t) = k) = \frac{(\lambda t)^k}{k!} e^{-\lambda t}, \text{ for } k = 0, 1, \dots$$

with  $\lambda$  the average number of events per time interval. Furthermore, in a Poisson process the number of events in disjunct time intervals are independent of each other.

Determine the autocorrelation function  $r_x(t_1, t_2)$ .

**Exercise 5.22** (3.2 from Statistical Digital Signal Processing and Modeling, Hayes)

Let  $x(n)$  be a stationary random process with zero mean and autocorrelation  $r_x(k)$ . We form the process  $y(n)$  as follows:

$$y(n) = x(n) + f(n)$$

where  $f(n)$  is a known deterministic sequence. Find the mean  $m_y(n)$  and the autocorrelation  $r_y(k, l)$  of the process  $y(n)$ .

**Exercise 5.23** (3.3 from Statistical Digital Signal Processing and Modeling, Hayes)

A discrete-time random process  $x(n)$  is generated as follows:

$$x(n) = \sum_{k=1}^p a(k)x(n-k) + w(n) \quad (5.18)$$

where  $w(n)$  is a white noise process with variance  $\sigma_w^2$ . Another process,  $z(n)$ , is formed by adding noise to  $x(n)$ ,

$$z(n) = x(n) + v(n) \quad (5.19)$$

where  $v(n)$  is white noise with a variance of  $\sigma_v^2$  that is uncorrelated with  $w(n)$ . You may assume however that  $\sigma_v = \sigma_w$ .

- (a) Find the power spectrum of  $x(n)$ .
- (b) Find the power spectrum of  $z(n)$ .

**Exercise 5.24** Let

- $y_1(n) = A \cos(n\omega_0)$  where  $A$  is a Gaussian random variable with mean  $m_A$  and variance  $\sigma_A^2$ .
- $y_2(n) = A \cos(n\omega_0 + \phi)$  where  $\phi$  is a random variable that is uniformly distributed between  $-\pi$  and  $\pi$ .
- $y_3(n) = u \cos(n) + v \sin(n)$  where  $u$  and  $v$  mutually independent random variables with mean  $m_u = m_v = 0$  and variance  $\sigma_u^2 = \sigma_v^2 = 1$ .
- $y_4(n) = u \sin(n) + v \cos(n)$  where  $u$  and  $v$  mutually independent random variables with mean  $m_u = m_v = 0$  and variance  $\sigma_u^2 = \sigma_v^2 = 1$ .

- (a) For both  $y_1(n)$  and  $y_2(n)$  determine whether the process is:
  - Wide-sense stationary.
  - Ergodic in its mean.
- (b) Are the processes  $y_3(n)$  and  $y_4(n)$  jointly wide-sense stationary?

## Chapter 6

# Filtering Stochastic Processes

**After studying this chapter you can:**

- determine under which condition the WSS property of a stochastic process is preserved by filtering that stochastic process by an LTI system
- determine the Cross-correlation function and the Cross spectrum between the input and output of an LTI system given its impulse response and the Auto-correlation function of the input.
- determine the Auto-correlation function and the Power spectrum of the output of an LTI system given its impulse response and the Auto-correlation function of the input.
- define a WSS white noise stochastic process
- compute using the Yule-Walker equations the Auto-correlation function for different (generally a finite number) lags of the output of an LTI system given as an ARMA, AR or MA model and the input assumed to be zero-mean white noise.

---

## 6.1 Introduction

Wide Sense Stationary (WSS) stochastic processes defined in Chapter 5 are described by their second (order) moment statistics. These can be their Auto- and Cross correlation function in the time domain or their Power- or Cross Spectrum in the frequency domain.

When considering the filtering of stochastic processes by a dynamic system, we are interested in this Chapter in how these second order moment statistics change by the filtering process.

Three parts are considered in this chapter. The first part focuses on the filtering by a general mixed causal, anti-causal LTI system. Then a tiny part defines a specific stochastic process namely zero-mean white noise (ZMWN). Many stochastic processes are derived from ZMWN by filtering this white noise signal by an LTI filter. In the third part we consider zero-mean white noise as an input to specific LTI systems. These LTI systems are governed by finite order difference equation and are the so-called ARMA (Auto-Regressive Moving Average) models. In addition to the definition of this model and how the Auto-correlation function and Power spectrum of the white noise filtered signal depends on the parameters of an ARMA model, two special variants of ARMA models are considered. The first being the so-called AR (Auto-regressive) model and the second the so-call MA (Moving Average) model.

The problems considered are of the type of so-called *forward modeling* nature. As we assume that the parameters of the LTI filter are given and the statistics (Auto-correlation function or its Power spectrum) of the input is also given.

The organisation of this chapter is as follows. Section 6.2 analyses the conditions on an LTI system such that the property of WSS is preserved by linear filtering a WSS stochastic process. It further derives expressions for the Auto-(Cross-) correlation functions of the output (between the output and input) of an LTI system when its input is a WSS stochastic process with given Auto-correlation function. The second concludes with providing the analogue results in the frequency domain. In Section 6.3 a special WSS stochastic process is defined, namely zero-mean white noise (ZMWN). Section 6.4 specializes the insights of Section 6.2 to the case the LTI system has a so-called ARMA, AR or MA (parametric) model structure and the input is ZMWN. Such parametric model structure model the input-output behavior via difference equation rather than a convolution between in general infinite length time sequences.

---

## 6.2 General mixed causal, anti-causal LTI systems

Consider the filtering of a stochastic process  $x(n)$  by an LTI system as depicted in Figure 6.1, with the transfer function  $H(z)$  given as:

$$H(z) = \sum_{\ell=-\infty}^{\infty} h(\ell)z^{-\ell} \quad h(\ell) \in \mathbb{C} \quad (6.1)$$

Unless otherwise specified we are considering zero-mean stochastic processes.



**Figure 6.1:** Schematic Representation of a LTI Discrete-Time system represented by the transfer function  $H(z)$  given as (6.1) and transforming its input  $x(n)$  to its output  $y(n)$ .

The question we are considering in this section is how the given statistics of a WSS stochastic process, i.e. its Auto-correlation function or its Power Spectrum change due to filtering. This question is first considered in the time domain in Section 6.2.1 and then in the frequency domain in Section 6.2.2.

A key question that precedes this analysis is under which conditions the filtered signal remains WSS. This is summarized in our first Theorem of this Chapter.

**Theorem 6.1** (WSS of LTI filtered signal). *Consider the filtering of the (input) stochastic process  $x(n)$  by an LTI system with transfer function  $H(z)$  given in (6.1), as depicted in Figure 6.1, then the output given by the convolution:*

$$y(n) = \sum_{\ell=-\infty}^{\infty} h(\ell)x(n-\ell) \quad (6.2)$$

is WSS provided that the following two conditions are satisfied:

1. The input signal  $x(n)$  is WSS, and,
2. The LTI system given by the transfer function  $H(z)$  is BIBO stable.

*Proof.* According to Definition 5.6, three things need to be proven. First we consider the mean of  $y(n)$ . This is given for the mean of the input denoted as  $\mu_x$  as:

$$\begin{aligned} E[y(n)] &= \sum_{\ell=-\infty}^{\infty} h(\ell)E[x(n-\ell)] = \sum_{\ell=-\infty}^{\infty} h(\ell)\mu_x \\ &= H(z)|_{z=1} \mu_x \end{aligned}$$

This value exists provided  $H(z)|_{z=1}$  converges. Since the system is BIBO stable, Lemma 2.6 shows that,

$$|H(z)|_{z=1}| = \left| \sum_{\ell=-\infty}^{\infty} h(\ell) \right| \leq \sum_{\ell=-\infty}^{\infty} |h(\ell)| < \infty$$

and the mean  $E[y(n)]$  exists and is constant. Since  $H(z)|_{z=1}$  exists (since it converges), it can be concluded that  $E[y(n)]$  is constant (not dependent on time  $n$ ) and equal to zero.

*Second,* we consider the Auto-correlation function  $r_y(n, n-k)$  and show that its argument does not dependent on the time index  $n$ . This is shown using the

fact that  $E[\cdot]$  is a linear operator and that the input  $x(n)$  is WSS by the following sequence of expressions,

$$\begin{aligned}
r_y(n, n-k) &= E[y(n)y^*(n-k)] \\
&= E\left[\sum_{\ell=-\infty}^{\infty} h(\ell)x(n-\ell) \sum_{p=-\infty}^{\infty} h^*(p)x^*(n-k-p)\right] \\
&= \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell)h^*(p)E[x(n-\ell)x^*(n-k-p)] \\
&= \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell)h^*(p)r_x(n-\ell, n-k-p) \\
&= \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell)h^*(p)r_x(k+p-\ell)
\end{aligned}$$

which does not depend on  $n$ .

Finally, we have to show that the variance of  $y(n)$  is finite. This variance is equal to:

$$|r_y(0)| = \left| \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell)h^*(p)r_x(p-\ell) \right|$$

By Property 5.10, the variance is bounded as,

$$\begin{aligned}
|r_y(0)| &\leq \left| \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell)h^*(p) \right| |r_x(0)| \\
&\leq \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} |h(\ell)||h(p)| |r_x(0)| \\
&< \infty
\end{aligned}$$

Where the last strict inequality follows from the BIBO stability of the system.  $\square$

### 6.2.1 The Auto- and Cross correlation Function after filtering

We consider the filtering of a WSS stochastic process  $x(n)$  as depicted in Figure 6.1 with the transfer function given as (6.1). We first consider deriving an expression for the Cross-correlation function between the output and the input of this filter. This is given in the following Theorem.

**Theorem 6.2** (Cross-correlation function of LTI filtered signal). *Consider the filtering of the (input) stochastic process  $x(n)$  by an LTI system with transfer function  $H(z)$  given in (6.1), as depicted in Figure 6.1, with the output given by the convolution in (6.2), then, the Cross-correlation  $r_{yx}(k)$  equals,*

$$\begin{aligned}
r_{yx}(k) &= \sum_{\ell=-\infty}^{\infty} h(\ell)r_x(k-\ell) \\
&= h(k) \star r_x(k)
\end{aligned} \tag{6.3}$$



and the Cross-correlation  $r_{xy}(k)$  equals,

$$\begin{aligned} r_{xy}(k) &= \sum_{\ell=-\infty}^{\infty} h^*(\ell) r_x(k + \ell) \\ &= h^*(-k) \star r_x(k) \end{aligned} \quad (6.4)$$

*Proof.* This is left as an exercise. See Exercise 6.1.  $\square$

Pay attention to the order of the indices of the Cross-correlation function. In (6.3) we have in the index the output variable  $y$  first and then the input variable  $u$ . This order of indices is reversed in (6.4).

The expression for the Auto-correlation function  $r_y(k)$  of the output signal is given in the following Theorem.

**Theorem 6.3** (Auto-correlation function of LTI filtered signal). *Consider the filtering of the (input) stochastic process  $x(n)$  by an LTI system with transfer function  $H(z)$  given in (6.1), as depicted in Figure 6.1, and the output given by the convolution in (6.2), then,*

$$\begin{aligned} r_y(k) &= \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell) h^*(p) r_x(k + p - \ell) \\ &= h(k) \star h^*(-k) \star r_x(k) \end{aligned} \quad (6.5)$$

*Proof.* The first part of the expression for  $r_y(k)$  is already given in the proof of Theorem 6.1. Based on this expression and the definition of the convolution as expressed in (2.11), we have that,

$$\begin{aligned} r_y(k) &= \sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell) h^*(p) r_x(k + p - \ell) \\ &= \sum_{\ell=-\infty}^{\infty} h(\ell) [h^*(\ell - k) \star r_x(k - \ell)] \\ &= h(k) \star h^*(-k) \star r_x(k) \end{aligned}$$

$\square$

## 6.2.2 The Power and Cross Spectrum after filtering

Again we first derive an expression for the Cross spectrum in Theorem 6.4.

**Theorem 6.4** (Cross Spectrum of LTI filtered signal). *Consider the filtering of the (input) stochastic process  $x(n)$  by an LTI system with transfer function  $H(z)$  given in (6.1), as depicted in Figure 6.1, and the output given by the convolution in (6.2), then, the Cross-Spectrum  $P_{yx}(z)$  is given as:*

$$P_{yx}(z) = H(z) P_x(z) \quad (6.6)$$

and  $P_{xy}(z)$  as:

$$P_{xy}(z) = H^*(1/z^*) P_x(z) \quad (6.7)$$

*Proof.* The proof of (6.6) directly follows from the z-transform of the convolution given in Table 2.3. The expression (6.7) is found by the following sequence of operations,

$$\begin{aligned}
P_{xy}(z) &= \sum_{k=-\infty}^{\infty} r_{xy}(k)z^{-k} \\
&= \sum_{k=-\infty}^{\infty} E[x(n)y^*(n-k)]z^{-k} \\
&= \sum_{k=-\infty}^{\infty} E[x(n+k)y^*(n)]z^{-k} \\
&= \sum_{k=-\infty}^{\infty} E[x(n+k) \sum_{\ell=-\infty}^{\infty} h^*(\ell)x^*(n-\ell)]z^{-k} \\
&= \sum_{\ell=-\infty}^{\infty} h^*(\ell) \left[ \sum_{k=-\infty}^{\infty} r_x(k+\ell)z^{-k} \right] \\
&= \sum_{\ell=-\infty}^{\infty} h^*(\ell)z^{\ell} P_x(z) \\
&= H^*(1/z^*)P_x(z)
\end{aligned}$$

□

The expression for the Power spectrum is given next.

**Theorem 6.5** (Power Spectrum of LTI filtered signal). *Consider the filtering of the (input) stochastic process  $x(n)$  by an LTI system with transfer function  $H(z)$  given in (6.1), as depicted in Figure 6.1, and the output given by the convolution in (6.2), then, the Power Spectrum  $P_y(z)$  is given as:*

$$P_y(z) = H(z)H^*(1/z^*)P_x(z) \quad (6.8)$$

*Proof.* The proof is left as an exercise (see Exercise 6.2). □

An overview of the results of Theorems 6.2, 6.3, 6.4 and 6.5 is given in Table 6.1.

---

## 6.3 Zero-Mean White Noise

An important way to generate stochastic processes, or viewed from a practical perspective, an important way to assume stochastic processes are generated is that they are a filtered version of some “standardized” stochastic process. Such a “standard” is often taken as Zero-Mean White Noise.

**Definition 6.6** (Zero-mean White Noise (ZMWN)). *A stochastic process  $x(n) \in \mathbb{C}$  is Zero-mean White Noise (ZMWN) provided it (1) has mean zero, (2) it is WSS and (3) its Auto-correlation function is given as,*

$$r_x(k) = \sigma_x^2 \delta(k)$$

| Quantity                      | Expression   |
|-------------------------------|--|
| Cross-correlation $r_{yx}(k)$ | $\sum_{\ell=-\infty}^{\infty} h(\ell)r_x(k-\ell)$<br>$= h(k) \star r_x(k)$   |
| Cross-correlation $r_{xy}(k)$ | $\sum_{\ell=-\infty}^{\infty} h^*(\ell)r_x(k+\ell)$<br>$= h^*(-k) \star r_x(k)$  |
| Auto-correlation $r_y(k)$     | $\sum_{\ell=-\infty}^{\infty} \sum_{p=-\infty}^{\infty} h(\ell)h^*(p)r_x(k+p-\ell)$<br>$= h(k) \star h^*(-k) \star r_x(k)$ |
| Cross Spectrum $P_{yx}(z)$    | $H(z)P_x(z)$   |
| Cross Spectrum $P_{xy}(z)$    | $H^*(1/z^*)P_x(z)$   |
| Auto Spectrum $P_y(z)$        | $H(z)H^*(1/z^*)P_x(z)$   |

**Table 6.1:** Summary of the Cross- and Autocorrelation functions and the Cross- and Power Spectra corresponding the LTI filtering considered in Figure 6.1 with the signal  $y(n)$  given by (6.2).

with  $\delta(k)$  defined in Table 2.1.

When we consider the ZMWN  $x(n)$  to be complex it may be specified as,

$$x(n) = x_{\text{Re}}(n) + jx_{\text{Im}}(n)$$

with  $x_{\text{Re}}(n)$  and  $x_{\text{Im}}(n)$  both real ZMWN and uncorrelated. Let their corresponding Auto-correlation functions be denoted as  $r_{x_{\text{Re}}}(k)$  resp.  $r_{x_{\text{Im}}}(k)$ , then,

$$\begin{aligned}
r_x(k) &= E[x(n)x^*(n-k)] \\
&= E\left[\left(x_{\text{Re}}(n) + jx_{\text{Im}}(n)\right)\left(x_{\text{Re}}(n-k) - jx_{\text{Im}}(n-k)\right)\right] \\
&= r_{x_{\text{Re}}}(k) + r_{x_{\text{Im}}}(k)
\end{aligned}$$

It is important to note that the definition of ZMWN does not specify the probability density function (pdf) of the individual samples for a fixed time instance. The WSS however requires these pdfs to be identical. For example at each time instance the samples could be a uniform random variable with a pdf as depicted in Figure 3.2 (on the right). Or it could be a Bernoulli random variable. The latter has been defined in Definition 3.2. An example of a ZMWN realization with the individual samples taken as Bernoulli random variables is plotted in Figure 5.1.

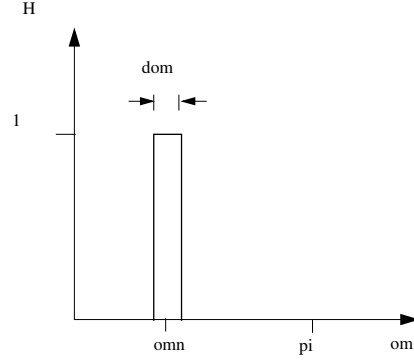
The Power spectrum of a ZMWN signal (or stochastic process) can be derived from the z-transform of the unit pulse given in Table 2.4. When  $x(n)$  is ZMWN with variance  $\sigma_x^2$  its Power Spectrum is:

$$P_x(e^{j\omega}) = \sigma_x^2$$

As the DTFT is a periodic function, it is completely defined on the interval  $[-\pi, \pi]$ . As such that Power spectrum is constant with value  $\sigma_x^2$  on that interval.

From the remark made in 5.22, it can be concluded that the Power of a ZMWN signal is equally distributed over all frequencies. That is when filtering the

ZMWN signal with an identical (and ideal) bandpass filter such that its transfer function  $H(z)$  has the magnitude of the Bode plot as depicted in Figure 6.2 for a center frequency  $\omega_0$ , then for each  $\omega_0 \in [\frac{\Delta\omega}{2}, \pi - \frac{\Delta\omega}{2}]$ , the filtered signal has identical Power.



**Figure 6.2:** The magnitude plot of an ideal Band-pass filter  $H(e^{j\omega})$ .

---

## 6.4 ARMA, AR, MA models

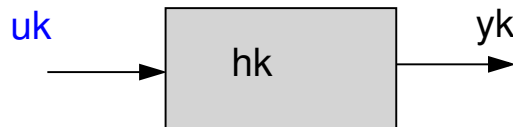
In physics but also in many other fields of engineering analysing stochastic processes, filtering is generally not done as described by the convolution in 6.2. As illustrated with the discretization of the Langevin equation as given by (1.10), finite difference equations are more frequently used to simulate dynamical systems. Here we could interpret the simulation actual as a way to generate stochastic processes. In this section 3 different widely used finite difference equations are discussed to generate stochastic processes by filtering ZMWN.

Three different variants of the so-called ARMA (Auto-Regressive Moving Average) analysed. ARMA models have first been described in [1] and early books on this topic include [2].

### 6.4.1 Definition of the models

#### ARMA model

Consider again the general (linear) filter schematic as depicted in Figure 6.3. Contrary to the general scenario considered in Figure 6.1, the input  $v(n)$  is re-



**Figure 6.3:** Schematic Representation of the input-output relationship in the definition of ARMA, AR and MA models.

stricted to ZMWN with variance  $\sigma_x^2$  and further the filter is assumed to have the

transfer function  $H(z)$  given as:

$$H(z) = \frac{\sum_{k=0}^q b(k)z^{-k}}{1 + \sum_{k=1}^p a(k)z^{-k}} \quad p \geq 1, q \geq 0 \quad (6.9)$$

This transfer function defines the transfer function of an ARMA model. The model is determined by the two polynomials  $A(z) = 1 + \sum_{k=1}^p a(k)z^{-k}$  and  $B(z) = \sum_{k=0}^q b(k)z^{-k}$ . These polynomials are of *finite order*  $p$  and  $q$  resp. and characterized by the parameters  $a(1), \dots, a(p)$  and  $b(0), \dots, b(q)$ .

By the definition of the transfer function (6.9), the relation between the time sequences  $v(n)$  and  $x(n)$  is given by the following difference equation:

$$x(n) + a(1)x(n-1) + \dots + a(p)x(n-p) = b(0)v(n) + b(1)v(n-1) + \dots + b(q)v(n-q) \quad (6.10)$$

**Lemma 6.7** (Taylor series and ARMA transfer function coefficients). *For the ARMA model the coefficients  $h(n)$  of the transfer function  $H(z)$  can be found using Taylor series:*

$$h(n) = \frac{H^{(n)}|_{z^{-1}=0}}{n!}$$

$$\text{where } H^{(n)}|_{z^{-1}=0} = \frac{\partial^n H(z)}{\partial (z^{-1})^n} \Big|_{z^{-1}=0}$$

*Proof.* All ARMA models are causal:

$$H(z) = \sum_{n=0}^{\infty} h(n)z^{-n}$$

Using the substitution  $x = z^{-1}$  we have:

$$\tilde{H}(x) = H(x^{-1}) = H(z)$$

And applying the Taylor series to the  $\tilde{H}(x)$ :

$$\tilde{H}(x) = \sum_{n=0}^{\infty} \frac{\tilde{H}^{(n)}(0)}{n!} x^n$$

Substituting back  $z^{-1} = x$ :

$$H(z) = \sum_{n=0}^{\infty} \frac{H^{(n)}|_{z^{-1}=0}}{n!} z^{-n}$$

$$\text{where } H^{(n)}|_{z^{-1}=0} = \frac{\partial^n H(z)}{\partial (z^{-1})^n} \Big|_{z^{-1}=0}$$

□

Another way to find  $h(n)$  is to use the Taylor series for fractions  $\frac{1}{1-x}$  as it will be shown later in the example 6.2.

### AR model

For the AR model we consider the same input-output set up as depicted in Figure 6.3 now with the only difference that the transfer function  $H(z)$  is given as:

$$H(z) = \frac{b(0)}{1 + \sum_{k=1}^p a(k)z^{-k}} \quad p \geq 1 \quad (6.11)$$

This transfer function defines the transfer function of an AR model. The model is determined by the one polynomial  $A(z) = 1 + \sum_{k=1}^p a(k)z^{-k}$ . The parameters that define an AR model are the *finite order*  $p$  of  $A(z)$ , its parameters  $a(1), \dots, a(p)$  and the parameter  $b(0)$ .

By the definition of the transfer function (6.11), the relation between the time sequences  $v(n)$  and  $x(n)$  is given by the following difference equation:

$$x(n) + a(1)x(n-1) + \dots + a(p)x(n-p) = b(0)v(n)$$

It is remarked that the discretization of the Langevin equation belongs to the class of AR models of order  $p = 2$ . These AR models are used abundantly in various engineering fields, but also in many other fields such as in econometrics [3].

### MA model

For the MA model we consider the same input-output set up as depicted in Figure 6.3 now with the only difference that the transfer function  $H(z)$  is given as:

$$H(z) = \sum_{k=0}^q b(k)z^{-k} \quad q \geq 0 \quad (6.12)$$

This transfer function defines the transfer function of an MA model. The model is determined by only one polynomial  $B(z) = \sum_{k=0}^q b(k)z^{-k}$ . The parameters that define an MA model are the *finite order*  $q$  of  $B(z)$  and its parameters  $b(0), \dots, b(q)$ .

By the definition of the transfer function (6.12), the relation between the time sequences  $v(n)$  and  $x(n)$  is given by the following difference equation:

$$x(n) = b(0)v(n) + b(1)v(n-1) + \dots + b(q)v(n-q)$$

## 6.4.2 Calculation of the Power Spectrum

For the 3 specific models we now proceed to compute the Power spectra of the output. Hereby making use of the assumption that the input signal  $v(n)$  as depicted in Figure 6.3 is ZMWN with variance  $\sigma_v^2$ .

### For the ARMA model

Consider the transfer function of the ARMA model be given as in (6.9), now denoted as:

$$H(z) = \frac{B_q(z)}{A_p(z)} \quad (6.13)$$

for  $A_p(z) = 1 + \sum_{k=1}^p a(k)z^{-k}$  and  $B_q(z) = \sum_{k=0}^q b(k)z^{-k}$ . Then applying the expression for the power spectrum for general transfer function  $H(z)$  as given in (6.8) of Theorem 6.5, we obtain the following expression for the Power Spectrum of the filtered output of an ARMA model:

$$P_x(z) = \frac{B_q(z)B_q^*(1/z^*)}{A_p(z)A_p^*(1/z^*)}\sigma_v^2 \quad (6.14)$$

For  $z = e^{j\omega}$  the Power Spectrum can be denoted as,

$$P_x(e^{j\omega}) = \frac{|B_q(e^{j\omega})|^2}{|A_p(e^{j\omega})|^2}\sigma_v^2 \quad (6.15)$$

It can indeed be checked, as requested for by Exercise 6.7 that the Power Spectrum as given by (6.14) and (6.15) for both the complex and real case (of signals and filter coefficients) satisfy Properties 5.19 and 5.20.

An illustration of a Power Spectrum for an ARMA model is given in the following example.

---

**Example 6.1 (Power Spectrum ARMA model)**

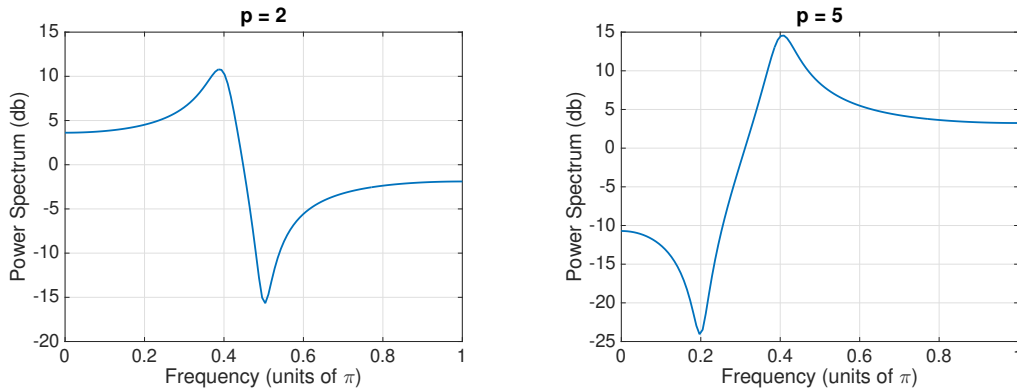
In this example we generate a 2nd order ARMA model with  $p = q = 2$ . For that purpose we define the complex numbers:

$$a = 0.95e^{j\frac{\pi}{f}} \quad b = 0.9e^{j\frac{2\pi}{5}} \quad \text{for } f \in \mathbb{N}$$

and let the transfer function be equal to,

$$H(z) = \frac{1 - (a + a^*)z^{-1} + |a|^2z^{-2}}{1 - (b + b^*)z^{-1} + |b|^2z^{-2}}$$

Then for  $f = 2$  and  $f = 5$  we display the Power Spectrum in Figure 6.4. It should



**Figure 6.4:** Power Spectrum of the ARMA model considered in Example 6.1 for 2 different values of  $f$ .

be observed from Figure 6.4 that for  $f = 2$  the Power Spectrum shows a peak at  $\omega = \frac{2\pi}{5}$ . This is a consequence of the pole location in  $0.95e^{j\frac{2\pi}{5}}$ . Its conjugate one

causes a peak at  $\omega = -\frac{2\pi}{5}$ . However since the Power Spectrum is symmetric this is not displayed here. Further for that value of  $f = 2$ , we observe that the Power Spectrum has a valley at  $\omega = \frac{\pi}{2}$ . This is caused by the location of the zero at  $0.95e^{j\frac{\pi}{2}}$ . For  $f = 5$  we see that the order of the peak and valley have changed.

From this example it may be concluded that for resonant ARMA models, i.e. ARMA models with poles and zeros close the unit circle, the Power Spectrum provides information on the pole and zero locations of the ARMA model.

---

### For the AR model

Consider the transfer function of the AR model be given as in (6.11), now denoted as:

$$H(z) = \frac{1}{A_p(z)}$$

for  $A_p(z) = 1 + \sum_{k=1}^p a(k)z^{-k}$ . Then applying the expression for the power spectrum for general transfer function  $H(z)$  as given in (6.8) of Theorem 6.5, we obtain the following expression for the Power Spectrum of the filtered output of an AR model:

$$P_x(z) = \frac{1}{A_p(z)A_p^*(1/z^*)}\sigma_v^2 \quad (6.16)$$

For  $z = e^{j\omega}$  the Power Spectrum can be denoted as,

$$P_x(e^{j\omega}) = \frac{1}{|A_p(e^{j\omega})|^2}\sigma_v^2 \quad (6.17)$$

One could ask the question why the parameter  $b(0)$  has been taken equal to one in the AR model in (6.16) for the expression for the Power Spectrum. It is clear from (6.17) that a parameter  $b(0)$  different by one could be interpreted as a change of the variance of the ZMWN.

### For the MA model

Consider the transfer function of the MA model be given as in (6.12), now denoted as:

$$H(z) = B_q(z)$$

for  $B_q(z) = \sum_{k=0}^q b(k)z^{-k}$ . Then applying the expression for the power spectrum for general transfer function  $H(z)$  as given in (6.8) of Theorem 6.5, we obtain the following expression for the Power Spectrum of the filtered output of an MA model:

$$P_x(z) = B_q(z)B_q^*(1/z^*)\sigma_v^2 \quad (6.18)$$

For  $z = e^{j\omega}$  the Power Spectrum can be denoted as,

$$P_x(e^{j\omega}) = |B_q(e^{j\omega})|^2\sigma_v^2 \quad (6.19)$$



### 6.4.3 Calculation of the Auto-Correlation Function

For the three different models (ARMA, AR and MA) defined in Section 6.4.2, we compute the Auto-Correlation Function of the output. This is done again for variance of the ZMWN input signal  $v(n)$  denoted by  $\sigma_v^2$ .

#### For the ARMA model

The equation from which the coefficients of the Auto-Correlation Function of the output of an ARMA model can be computed is called the *Yule-Walker* equation. This famous equation is listed in the following Theorem.

**Theorem 6.8** (Yule-Walker equation). *Let the WSS stochastic process  $x(n) \in \mathbb{C}$  be generated by filtering ZMWN with variance  $\sigma_v^2$  as depicted in Figure 6.3 with the transfer function  $H(z)$  defined as for the ARMA model (6.9), with possibly complex coefficients, let  $H(z)$  be BIBO stable and let the transfer function  $H(z)$  also be expanded as,*

$$H(z) = \sum_{k=0}^{\infty} h(k)z^{-k} \quad (6.20)$$

then the Auto-correlation function  $r_x(k)$  satisfies the following Yule-Walker equation:

$$r_x(n) + \sum_{\ell=1}^p a(\ell)r_x(n-\ell) = \begin{cases} \sigma_v^2 \sum_{\ell=n}^q b(\ell)h^*(\ell-n) & : 0 \leq n \leq q \\ 0 & : n > q \end{cases} \quad (6.21)$$

*Proof.* Since  $H(z)$  is assumed to be BIBO stable, the output  $x(n)$  is WSS (see Theorem 6.1). As such its Auto-correlation function can be denoted as  $r_x(k)$ . To find the expression (6.21) we start however from its Power Spectrum.

The Power spectrum of the output  $x(n)$  as depicted in Figure 6.3 is for general  $H(z)$  (and hence also for ARMA models) given as in (6.8). Using the ZMWN property of the input  $v(n)$ , the Power Spectrum becomes:

$$P_x(z) = \sigma_v^2 H(z)H^*(1/z^*) = \sigma_v^2 H^*(1/z^*)H(z) \quad (6.22)$$

For the ARMA model  $H^*(1/z^*)$  equals:

$$H^*(1/z^*) = \frac{b(0)^* + b(1)^*z + \cdots + b(q)^*z^q}{1 + a(1)^*z + \cdots + a(p)^*z^p}$$

Based on this expression of  $H^*(1/z^*)$  we can write (6.22) as,

$$P_x(z) + a(1)^*zP_x(z) + \cdots + a(p)^*z^pP_x(z) = \sigma_v^2 \left( b(0)^*H(z) + b(1)^*zH(z) + \cdots + b(q)^*z^qH(z) \right)$$

Taking the inverse of the z-transform as indicated in Remark 2.1, the Auto-correlation satisfies the following difference equation,

$$r_x(k) + a(1)^*r_x(k+1) + \cdots + a(p)^*r_x(k+p) = \sigma_v^2 \left( b(0)^*h(k) + b(1)^*h(k+1) + \cdots + b(q)^*h(k+q) \right)$$

If we take the complex conjugate of both sides, we obtain,

$$r_x^*(k) + a(1)r_x^*(k+1) + \cdots + a(p)r_x^*(k+p) = \sigma_v^2 \left( b(0)h^*(k) + b(1)h^*(k+1) + \cdots + b(q)h^*(k+q) \right)$$

Further since,  $r_x^*(k) = r_x(-k)$ , we can write this equation as,

$$r_x(-k) + a(1)r_x(-k-1) + \cdots + a(p)r_x(-k-p) = \sigma_v^2 \left( b(0)h^*(k) + b(1)h^*(k+1) + \cdots + b(q)h^*(k+q) \right)$$

Now change the argument  $-k$  by  $n$  and we obtain,

$$r_x(n) + a(1)r_x(n-1) + \cdots + a(p)r_x(n-p) = \sigma_v^2 \left( b(0)h^*(-n) + b(1)h^*(-n+1) + \cdots + b(q)h^*(-n+q) \right)$$

This can compactly be denoted as,

$$r_x(n) + \sum_{\ell=1}^p a(\ell)r_x(n-\ell) = \sigma_v^2 \sum_{\ell=0}^q b(\ell)h^*(\ell-n)$$

Since the ARMA model is causal and BIBO,  $H(z)$  can be expanded as given in the theorem. This indicates that  $h(\ell-n) = 0$  for  $\ell < n$ . Taking into account the fact that for  $n > q$   $b(n) = 0$ , the last equation can be written as the Yule-Walker equation (6.21).  $\square$

**Remark 6.9.** *It should be remarked that using the Symmetry property 5.8 of the Auto-correlation function it is only necessary to find the samples  $r_x(n)$  for  $n \geq 0$ .*

For one value of  $n$  we obtain an equation with  $p+1$  unknowns. Inspecting the Yule-Walker equation we can take  $n = 0, 1, \dots, p$  and arrive at a set of  $p+1$  equations with  $p+1$  unknown. The latter unknowns being the values of the Auto-correlation function  $r_x(0), r_x(1), \dots, r_x(p)$ . The solution of these equations serves as the initial conditions for a recursion for  $r_x(n)$  that is given by the Yule-Walker equation for  $n > p$ . From this recursion and the symmetry property of the Auto-correlation function all other values of that function are determined.

The use of Theorem 6.8 is illustrated in the following example. A more complex application is called for in Exercise 6.9

---

### Example 6.2 (Auto-correlation function of ARMA model)

Consider the ARMA model with transfer function as given in Exercise 6.7 and repeated here:

$$H(z) = \frac{1 - 0.5z^{-1}}{1 - 0.8z^{-1}}$$

Let us consider  $\sigma_v = 1$ . Then the expansion for  $H(z)$  given as in (6.20), can be found by using Euclid's algorithm. This yields,

$$\frac{1 - 0.5z^{-1}}{1 - 0.8z^{-1}} = 1 + 0.3z^{-1} + 0.24z^{-2} + \cdots$$

Another way to find the  $h(n)$  is to use the Taylor expansion:

$$\begin{aligned}
H(z) &= (1 - 0.5z^{-1}) \sum_{n=0}^{\infty} (0.8)^n z^{-n} \\
&= \sum_{n=0}^{\infty} (0.8)^n z^{-n} - \sum_{n=0}^{\infty} 0.5(0.8)^n z^{-n-1} \\
&= \sum_{n=0}^{\infty} (0.8)^n z^{-n} - \sum_{n=1}^{\infty} 0.5(0.8)^{n-1} z^{-n} \\
&= 1 + \sum_{n=1}^{\infty} \frac{3}{8} (0.8)^n z^{-n} \\
&= 1 + 0.3z^{-1} + 0.24z^{-2} + \dots
\end{aligned}$$

Therefore  $h(0) = 1$ ,  $h(1) = 0.3$  and the Yule-Walker equation for  $n = 0$  reads:

$$r_x(0) - 0.8r_x(-1) = h(0) - 0.5h(1)$$

Using the symmetry of the Auto-correlation function  $r_x(k)$ , we can denote this as:

$$r_x(0) - 0.8r_x(1) = h(0) - 0.5h(1)$$

here we used the fact that  $h(n)$  is a real numbers, auto-correlation function of ZMWN  $r_v(k)$  is a real function, so from table 6.1 the result  $r_x(k) = h(k) \star h^*(-k) \star r_v(k)$  also is a real function.

For  $n = 1$  the Yule-Walker equation becomes:

$$r_x(1) - 0.8r_x(0) = 0 - 0.5h(0)$$

The last 2 equations for the computed values of  $h(0)$  and  $h(1)$  can be denoted in the matrix format,

$$\begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix} \begin{bmatrix} r_x(0) \\ r_x(1) \end{bmatrix} = \begin{bmatrix} 1 - .15 \\ -0.5 \end{bmatrix}$$

From which it follows that  $r_x(0) = \frac{5}{4}$  and  $r_x(1) = \frac{1}{2}$ . The other samples of  $r_x(k)$  result from the Yule-Walker equation for  $n \geq 2$ , being the following recursion:

$$r_x(n) - 0.8r_x(n-1) = 0 \quad n \geq 2$$

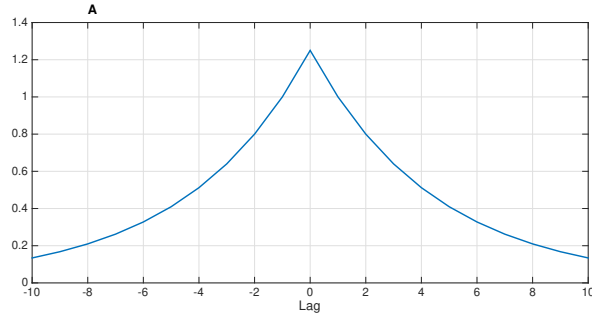
This yields the samples  $r_x(n)$  for  $n = -10$  to  $n = 10$  as plotted in Figure 6.5.

---

### For the AR model

As the AR model is a special case of the ARMA model with the following parameters:

$$q = 0 \quad \text{and} \quad b(0) = 1$$



**Figure 6.5:** Auto-correlation function computed with the Yule-Walker equation in Example 6.2.

the first parameter  $h(0)$  of the expansion as given in (6.20) equals 1 and therefore the Yule-Walker equations (6.21) become:

$$r_x(n) + \sum_{\ell=1}^p a(\ell)r_x(n-\ell) = \sigma_v^2\delta(n) \quad n \geq 0 \quad (6.23)$$

As for the ARMA model this equation is again used to build a set of  $p+1$  equations in the  $p+1$  unknowns of the samples of the Auto-correlation function  $r_x(n)$  for  $n = 0, 1, \dots, p$ . Then for  $n > p$  we use the recursion,

$$r_x(n) + \sum_{\ell=1}^p a(\ell)r_x(n-\ell) = 0, \quad n > p$$

and the found values of the Auto-correlation function to find (one-side) of the Auto-correlation function.

### For the MA model

Also the MA model is a special case of the ARMA model with the following parameters:

$$p = 0$$

For this case the parameters  $h(k)$  in the expansion (6.20) simply are,

$$h(k) = b(k) \quad 0 \leq k \leq q$$

Hence, the Yule-Walker equations (6.21) become:

$$r_x(n) = \begin{cases} \sigma_v^2 \sum_{\ell=n}^q b(\ell)b^*(\ell-n) & : 0 \leq n \leq q \\ 0 & : n > q \end{cases} \quad (6.24)$$

The Yule-Walker equation (6.24) provides values for the Auto-correlation function  $r_x(n)$  that are different from zero.

---

## References

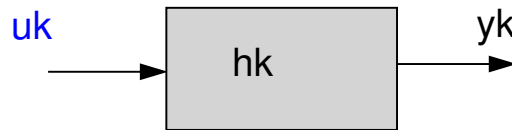
- [1] P. Whittle, *Hypothesis testing in time series analysis*. Uppsala, Almqvist & Wiksells boktr, 1951
  - [2] G. Box, G.M. Jenkins and G.C. Reinsel. *Time Series Analysis: Forecasting and Control*, (Third ed.). Prentice-Hall., 1994.
  - [3] E.J. Hannan and M. Deistler, *The statistical theory of linear systems*, Wiley, 1988.
- 

## Exercises

**Exercise 6.1** Prove Theorem 6.2.

**Exercise 6.2** Prove Theorem 6.5.

**Exercise 6.3** When  $y(n)$  is generated by filtering a ZMWN signal  $x(n)$  with unit variance ( $\sigma_x^2 = 1$ ) and when we have the Power Spectrum of the output  $y(n)$  given, show how to derive the magnitude of the Bode plot of  $G(e^{j\omega})$ .



**Figure 6.6:** The filter  $G(z)$  to generate  $y(n)$  in Exercise 6.3 and 6.4.

**Exercise 6.4** When  $y(n)$  is generated by filtering a ZMWN signal  $x(n)$  with unit variance ( $\sigma_x^2 = 1$ ) and when we have the Cross Spectrum between the output  $y(n)$  and the input  $x(n)$  given as  $R_{xy}(e^{j\omega})$ , then show how to derive  $G(e^{j\omega})$ .

**Exercise 6.5** Let the Cross Spectrum between the signals  $y(n)$  and  $x(n)$  be given as  $P_{yx}(z)$ . Let the signal  $x(n)$  be filtered by the filter with transfer function  $G(z)$  and call the output  $v(n)$ , such that:

$$V(z) = G(z)X(z)$$

then determine the Cross Spectrum  $P_{yv}(z)$  between the signals  $y(n)$  and  $v(n)$ . Here you may not assume that the filter  $G(z)$  is invertible.

**Exercise 6.6** Let the Cross Spectrum between the signals  $y(n)$  and  $x(n)$  be given as  $P_{yx}(z)$ . Let the signal  $y(n)$  be filtered by the filter with transfer function  $G(z)$  and call the output  $u(n)$ , such that:

$$U(z) = G(z)Y(z)$$

then determine the Cross Spectrum  $P_{ux}(z)$  between the signals  $u(n)$  and  $x(n)$ .

**Exercise 6.7** Consider the ARMA model with input ZMWN with variance  $\sigma_v$  and with transfer function given as:

$$H(z) = \frac{1 - 0.5z^{-1}}{1 - 0.8z^{-1}}$$

then

- (a) check whether the output of this model is WSS.
- (b) compute for  $\sigma_v = 1$  the Power Spectrum as given by (6.14) and (6.15) for this specific ARMA model.
- (c) Show that the Power Spectrum derived from (6.15) equals:

$$P_x(e^{j\omega}) = \frac{1.25 - \cos \omega}{1.64 - 1.6 \cos \omega}$$

Use this result to show that  $P_x(e^{j\omega})$  indeed satisfies Properties 5.19 and 5.20.

**Exercise 6.8** Let the AR model be given by the transfer function:

$$H(z) = \frac{1}{1 - az^{-1}} \quad a \in ]-1, 1[$$

then show that the Power Spectrum of the output of this AR model equals:

$$P_x(e^{j\omega}) = \frac{\sigma_v^2}{1 + a^2 - 2a \cos \omega}$$

**Exercise 6.9** Consider the ARMA model with transfer function given as:

$$H(z) = \frac{b(0) + b(1)z^{-1} + b(2)z^{-2} + b(3)z^{-3}}{1 + a(1)z^{-1} + a(2)z^{-2} + a(3)z^{-3}}$$

It may also be assumed that this transfer function is BIBO stable and that it can be written as,

$$H(z) = \sum_{m=0}^{\infty} h(m)z^{-m}$$

Assume that  $\sigma_v^2 = 1$ , then show that the Yule-Walker equations can be written as,

$$\begin{bmatrix} 1 & a(1) & a(2) & a(3) \\ a(1) & 1 + a(2) & a(3) & 0 \\ a(2) & a(1) + a(3) & 1 & 0 \\ a(3) & a(2) & a(1) & 1 \end{bmatrix} \begin{bmatrix} r_x(0) \\ r_x(1) \\ r_x(2) \\ r_x(3) \end{bmatrix} = \begin{bmatrix} \sum_{\ell=0}^3 b(\ell)h^*(\ell) \\ \sum_{\ell=1}^3 b(\ell)h^*(\ell-1) \\ \sum_{\ell=2}^3 b(\ell)h^*(\ell-2) \\ b(3)h^*(0) \end{bmatrix}$$

**Exercise 6.10** Consider the MA model to be given by the transfer function:

$$H(z) = b(0) + b(1)z^{-1} + b(2)z^{-2} \quad b(0), b(1), b(2) \in \mathbb{C}$$

consider  $\sigma_v^2 = 1$ , then show that the Auto-correlation function of the filtered signal  $x(n)$  is given by:

$$\begin{aligned} r_x(0) &= |b(0)|^2 + |b(1)|^2 + |b(2)|^2 \\ r_x(1) &= b(1)b^*(0) + b(2)b^*(1) \\ r_x(2) &= b(2)b^*(0) \end{aligned}$$

**Exercise 6.11** Let the relation between input  $u(t)$  and output  $y(t)$  of a causal discrete-time system be given by:

$$y(t) = u(t) - 0.5u(t-1) + 0.2u(t-2).$$

- (a) Compute the transfer function  $G(z)$  of the system.
- (b) Assume the input  $u(t)$  is a stationary stochastic process described by:

$$u(t) = 0.6e(t-1) + e(t), \quad (6.25)$$

where  $e(t)$  is a zero-mean white noise process with variance 1.

Compute the autocorrelation function  $r_u(\tau)$  of the input  $u(t)$ .

- (c) Compute the expected value of the output  $y(t)$ .
- (d) Compute the power spectrum  $P_y(\omega)$  of the output.

**Exercise 6.12** Let the relation between the input  $u(t)$  and the output  $y(t)$  of a linear causal discrete time system be described by:

$$y(t) = \sum_{k=0}^{\infty} g(k)u(t-k), \quad t = 0, 1, 2, \dots$$

with  $\{g(k)\}_{k=0,1,\dots}$  the impulse response of the system. Consider the case where

$$g(k) = \begin{cases} 1/2 & k = 0, 1 \\ 0 & \text{otherwise} \end{cases}$$

Suppose the input  $u(t)$  is a stationary stochastic process with expected value  $\mu_u = 0$  and autocorrelation function

$$R_u(\tau) = \begin{cases} 4 & \tau = 0, \\ 2 & \tau = \pm 1, \\ 0 & |\tau| \geq 2. \end{cases}$$

- (a) Compute the power spectrum  $P_u(\omega)$  of the input  $u(t)$ .
- (b) Can the expected value and the autocorrelation of the output  $y(t)$  be time-dependent?
- (c) Compute the expected value of the output  $y(t)$ .
- (d) Compute the autocorrelation function  $r_y(t, t-\tau)$  of the output  $y(t)$ .
- (e) Compute the variance of the output  $y(t)$ .

(f) Is the output  $y(t)$  wide sense stationary?

**Exercise 6.13** Consider a discrete time stationary stochastic process  $x(t)$  described by the following difference equation:

$$x(t) + a_1 x(t-1) = e(t) + b_1 e(t-1)$$

with  $a_1$  and  $b_1$  real constants and  $e(t)$  white noise with variance  $\sigma_e^2$ .

- (a) Compute  $r_x(2)$ , the autocorrelation of  $x$  for lag  $\tau = 2$ , in case  $a_1 = 0$ .
- (b) Compute  $r_x(2)$  in case  $b_1 = 0$ .
- (c) Compute the variance of the process  $x(t)$  (for arbitrary values of  $a_1$  and  $b_1$ ).

**Exercise 6.14** The relation between input  $u(t)$  and output  $y(t)$  of a linear causal discrete time system is given by:

$$y(t) = \sum_{k=0}^{\infty} g(k)u(t-k), \quad t = 0, 1, 2, \dots \quad (6.26)$$

with  $\{g(k)\}_{k=0,1,\dots}$  the impulse response of the system. Consider the case where

$$g(k) = \begin{cases} 1 & k = 0 \\ \frac{1}{2} & k = 1 \\ 0 & \text{otherwise} \end{cases}$$

Consider the case where  $u$  is a white noise process with variance 1.

- (a) Compute the cross-correlation function  $r_{uy}(\tau)$  between the input and output of the system.
- (b) Compute the power spectrum  $P_y(\omega)$  of the output of the system.
- (c) Now consider the case where the input  $u$  is a stationary stochastic process with expected value  $\mu_u = 0$  and power spectrum

$$P_u(\omega) = 1.04 + 0.2e^{i\omega} + 0.2e^{-i\omega}$$

Compute the autocorrelation function  $r_y(t, t-\tau)$  of the output  $y(t)$ .

- (d) Compute the cross spectrum  $P_{yu}(\omega)$  between the output and the input.

**Exercise 6.15** (3.11 from Statistical Digital Signal Processing and Modeling, Hayes)

Consider a first-order AR process that is generated by the difference equation

$$y(n) = ay(n-1) + w(n)$$

where  $|a| < 1$  and  $w(n)$  is a zero mean white noise random process with variance  $\sigma_w^2$ .



- (a) Give the unit sample response of the filter that generates  $y(n)$  from  $w(n)$ .
- (b) Give an expression for the autocorrelation function  $r_y(n)$  of  $y(n)$ .
- (c) Give an expression for the power spectrum  $P_y(\omega)$  of  $y(n)$ .

**Exercise 6.16** Consider an  $q$ th order moving average process that is generated by the difference equation:

$$y(n) = \sum_{k=0}^q b(k)w(n-k)$$

with  $w(n)$  a zero mean white noise with variance  $\sigma_w^2$ .

- (a) Give the unit sample response of the filter that generates  $y(n)$  from  $w(n)$ .
- (b) Give an expression for the autocorrelation function  $r_y(k)$  of  $y(k)$ .
- (c) Give an expression of the power spectrum  $P_y(e^{j\omega})$ .

**Exercise 6.17** Consider the random process:

$$x(n) = A \cos(n\omega_0 + \phi) + w(n)$$

where  $w(n)$  is zero mean white Gaussian noise with a variance  $\sigma_w^2$ . For each of the following cases, find an expression for the auto-correlation function, and if the process is WSS, find the power spectrum.

- (a)  $A$  is a Gaussian random variable with zero mean and variance  $\sigma_A^2$  and both  $\omega_0$  and  $\phi$  are constants.
- (b)  $\phi$  is uniformly distributed over the interval  $[-\pi, \pi]$  and both  $A$  and  $\omega_0$  are constants.
- (c)  $\omega_0$  is a random variable that is uniformly distributed over the interval  $[b - \Delta, b + \Delta]$  and both  $A$  and  $\phi$  are constants.



## Chapter 7

# Inverse Problems in Time and Frequency Domain

**After studying this chapter you can:**

- formulate and solve (a simple variant of) the Spectral Factorization problem
- formulate and solve the inverse problem to determine the coefficients of an ARMA resp. AR or MA model of a given order from an Auto-correlation function that is assumed to be the Auto-correlation function of the output of an LTI system with an ARMA resp. AR or MA model structure (but unknown) with ZMWN as input.

---

## 7.1 Introduction

In Chapter 6 we studied the filtering of stochastic processes. These type of problems belong to the field of so-called *forward modeling* problems or *simulation* of dynamical systems. We considered a model of a dynamical system in terms of the transfer function  $H(z)$  given in addition to the (statistical properties of the) input and we studied the question to compute the (statistical properties of the) output. The statistical properties considered were the Power Spectrum and Cross Spectrum in the Frequency domain and the Auto- and Cross-correlation function in the time domain. An interesting challenge in the context of stochastic processes is the generation of white noise as a discretization of its continuous counterpart. This is a theoretical challenging problem requiring the theory of so-called Ito Integrals [1] but for a more engineering oriented approach e.g explaining how to simulate the discrete Langevin equation (1.10) we refer to [3].

In this chapter we start with an introduction to the field of so-called *inverse problems*. Inverse problems refer to the field of science of calculating from a set of observations a (dynamic) mathematical model that can accurately reproduce these observations. Since the underlying physical mechanism of nature that has generated these observations can be very complex, one is generally interested in mathematical models of low complexity while approximating the observations accurately. The field of inverse problems starting from real-life measurements of physical processes is generally called *System Identification* and many good introductory books on this topic exist. See e.g. [4], [5] for more information on this topic.

In the context of this introduction a number of simple inverse problems are considered. The statistical information that will be assumed to be available and that constitutes the given observations are the Auto-correlation function of a stochastic process or its Power Spectrum. Based on this information the goal of the inverse problem is to find an ARMA model (see Section 6.4) and the variance of the white noise input, so that the Auto-correlation function of that model output or its Power Spectrum matches the given observation. All problems considered in this Chapter are restricted to the scalar case where a single Power spectrum or Auto-correlation function is given.

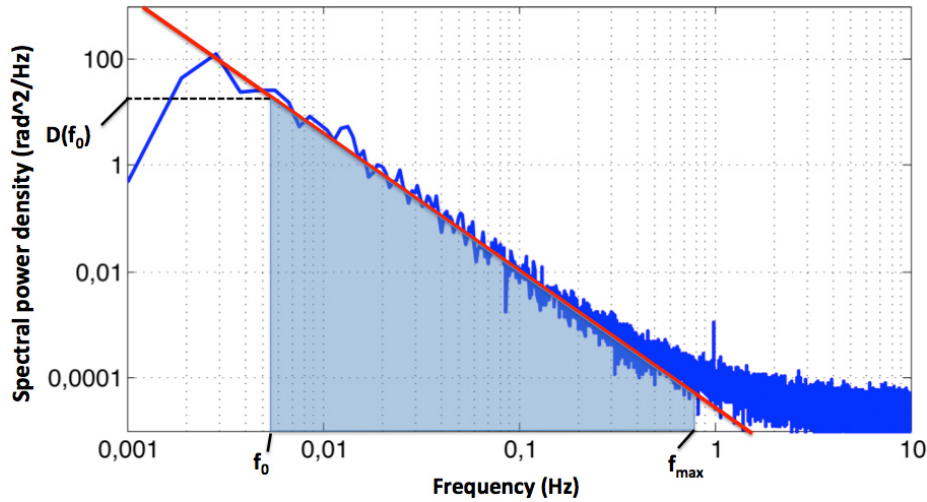
Another type of inverse problem is considered in the following Chapter 8 to estimate the parameters of an AR model from experimental (finite length) data sets. This is a simple introduction to the field of system identification.

Prior to the analysis of the two inverse problems in this Chapter we start with a motivating example in the Section 7.2. Then we treat the inverse problem based on the Power Spectrum in Section 7.3. The problem and its solution are part of the co-called class of Spectral Factorization problems. The inverse problem based on the Auto-correlation function is discussed in Section 7.4. In both Inverse problems the ARMA model that is to be determined will be indicated as the *Shaping Filter*.

## 7.2 A Motivating Example

Turbulence is an important physical phenomenon that occurs on a very wide scale. Examples include the flow around a golf ball, speed boat, wind turbines and even the heart sound picked up by a stethoscope is produced by turbulent blood flow. In optics turbulence may cause severe degradation of the image quality e.g. in ground based telescopes. A whole field known under the name Adaptive Optics exists in order to study this physical aberration and develop solutions to compensate for with it. An introduction to this field can e.g. be found in [6]. These are the course notes of the MSc course sc4045, Control for High Resolution Imaging.

A pioneering physicist in the field of turbulence modeling was the mathematician Andrej Kolmogorov. In his research to model turbulence he was the first to make use of Power Spectra measurements and to tried to fit a parametric model on these measurements in order to gain insight which parameters influence this physical phenomena. Typical data of such a Spectrum is displayed in Figure 7.1. The red curve in the figure is a fit to the Power spectrum in the



**Figure 7.1:** The Power Spectrum of phase scintillations  $\phi$  recorded in [7]. The blue “oscillatory” curve represents the real-life measurements and the straight line is a linear fit (in the log-log representation) that represents Kolmogorov’s  $\frac{5}{3}$  Power Law.

frequency range indicated by  $[f_0, f_{\max}]$ . Such fits were used by Kolmogorov to deduce his famous  $5/3$ -power law. Such law would indicate that the Power spectrum falls off as a function of the frequency  $f$  given by  $f^{-\frac{5}{3}}$ . When denoting such an approximation of the Power spectrum of the quantity of interest e.g. as  $P_\phi$ , it could then be denoted as:

$$P_\phi(f) = 0.023r_0^{-\frac{5}{3}}f^{-\frac{5}{3}} \quad (7.1)$$

with  $r_0 \in \mathbb{R}$  the so-called Fried Parameter, defined e.g. in [6], and expressing the strength of the turbulence.

The model developed by Kolmogorov was of interest to determine the strength and influence of turbulence in dynamics excitations to physical systems. As such it was important to know the distribution of the Power of turbulence over the frequency interval of interest in which the physical system is operational. Such experimental models are still of key interest and updated to actual flow scenarios in a wide variety of disciplines. To just name a few we have modeling aberrations due to the eye in retinal imaging or the turbulence-induced vibrations in wind turbines.

Though the parameters displayed in the above fit bears physical interest, we would like to go a step further in the analysis. We would like to artificially simulate a stochastic process in such a way that the simulated process has a Power spectrum that accurately approximates the experimentally observed one. This artificial stochastic process would allow to test new developments in simulation before doing (expensive) hardware testing. The topic fits in the framework of virtual prototyping.

As this is a problem of great complexity we study in this introductory book a simplified version. This is done in the next chapter starting from Power Spectrum data and in Section 7.4 from Auto-correlation data.

---

## 7.3 Spectral Factorization

### 7.3.1 Problem Definition

The starting point is a given (scalar) function  $F(z)$  in the complex variable  $z \in \mathbb{C}$  or for  $z = e^{j\omega}$  ( $\omega \in \mathbb{R}$ ) and the goal is to derive from this function a transfer function  $G(z)$  such that the following factorization of  $F(z)$  holds:

$$F(z) = s_0 G(z) G^*(1/z^*) \quad s_0 \in \mathbb{R}^+ \quad (7.2)$$

Recall that this factorization resembles the expression given for the Power Spectrum in Theorem 6.5. In order to find such a factorization we need to refine the conditions on the given function  $F(z)$  as well as have to stipulate more precise conditions on the transfer function  $G(z)$ . We start by stating the Assumptions the given function  $F(z)$  needs to satisfy for the simplified Spectral Factorization problem studied in this book.

### Assumptions

The Spectral Factorization problem considered here makes the following Assumptions on the given function  $F_x(z)$  or  $F(e^{j\omega})$ :

1.  $F_x(z)$  is a rational, coprime function in  $z$ , i.e. it can be written as the ratio of 2 finite order polynomials in  $z$  and these polynomials do not have common roots. Furthermore  $F_x(z)$  does not have poles on the unit circle.
2.  $F_x(z)$  is symmetric as stated in Property 5.19, repeated here:

$$F_x(z) = F_x^*(1/z^*) \quad (7.3)$$

3.  $F(e^{j\omega})$  is *positive real*, i.e.

$$F_x(e^{j\omega}) > 0 \quad \text{and} \quad F_x(e^{j\omega}) \in \mathbb{R} \quad (7.4)$$

This assumption is intrinsic to the definition of Power Spectra of WSS stochastic processes as highlighted in Property 5.20.

### Spectral Factorization Problem

For the above Assumptions on the given function  $F_x(z)$ , the *problem* is now to find the rational (transfer) function  $Q(z)$  that satisfies the following 4 conditions:

1.  $Q(z)$  is a *spectral factor* of  $F_x(z)$ , i.e.

$$F_x(z) = \sigma Q(z)Q^*(1/z^*) \quad (7.5)$$

2.  $Q(z)$  is a rational function that represent a causal and stable transfer function, i.e. its poles are within the unit circle.

3.  $Q(z)$  is minimum phase, i.e.  $Q(z)^{-1}$  is causal and stable, i.e. the zeros of  $Q(z)$  are within the unit circle.

4.  $Q(\infty) = 1$ .

The particular transfer function  $Q(z)$  will be called the *spectral factor*.

Remark that the first condition on the Spectral Factor  $Q(z)$  resembles the factorization given at the beginning of this section in (7.2). However as given in (7.2) the problem is not well-defined as it may e.g. give rise to many solutions. The additional requirements on  $Q(z)$  are meant to derive a *unique* solution.

After studying how to compute the spectral factor from the given function  $F_x(z)$  in the next subsection, we analyse its use in subsection 7.3.3.

### 7.3.2 Solution

In order to derive a solution we start exploring the consequences of the Assumptions stipulated on the function  $F_x(z)$ . The key assumption is the symmetry property (7.3). This in combination with  $F_x(z)$  to be rational imposes the coefficients of  $F_x(z)$  to be constrained. This is illustrated in the following example.

---

#### Example 7.1 (Constraints on the coefficients of $F_x(z)$ )

First consider  $F_x(z)$  to be a polynomial with complex coefficients given as:

$$F_x(z) = az + b + cz^{-1}$$

The function  $F_x^*(1/z^*)$  is also a polynomial given as:

$$F_x^*(1/z^*) = a^*z^{-1} + b^* + c^*z$$

By the symmetry property we have  $F_x(z) = F_x^*(1/z^*)$ . Then equating equal powers of  $z$  yields the following constraints on the coefficients:

$$a = c^* \quad b = b^*$$

An example of such a polynomial would be:

$$(1 + i)z + 4 + (1 - i)z^{-1}$$

As a second example, we consider the following rational function in  $z$  with real coefficients:

$$F_x(z) = \frac{az + b + cz^{-1}}{dz + e + fz^{-1}} \left( = \frac{a + bz^{-1} + cz^{-2}}{d + ez^{-1} + fz^{-2}} \right)$$

with  $a, b, c, d, e, f \in \mathbb{R}$ . For this rational function  $F_x^*(1/z^*)$  equals,

$$F_x^*(1/z^*) = \frac{az^{-1} + b + cz}{dz^{-1} + e + fz}$$

By the symmetry property (7.3), we get the following two relationships:

$$\begin{aligned} ae + bf &= bd + ce \\ af &= cd \end{aligned} \tag{7.6}$$

The last result yields for  $af \neq 0$  the following fractional relationship:

$$\frac{a}{d} = \frac{c}{f}$$

If we substitute this in the first relation of (7.6) we obtain, under the condition that  $a \neq c$  or  $d \neq f$ :

$$\frac{a}{d} = \frac{c}{f} = \frac{b}{e}$$

Exercise 7.1 calls for a derivation of this result.

---

More important than the constraints on the coefficients of the polynomials that define the given function  $F_x(z)$ , the symmetry condition (7.3) constrains the pole-zero location of such rational functions. This is summarized in the following Lemma.

**Lemma 7.1** (Conjugate reciprocal pole-zeros). *Let the complex function  $F_x(z) \in \mathbb{C}$  for  $z \in \mathbb{C}$  be rational and coprime, that satisfies,*

$$F_x(z) = F_x^*(1/z^*)$$

*then the following 2 statements hold,*

$$\begin{aligned} \text{If } p_0 \text{ is a pole of } F_x(z) &\Rightarrow \frac{1}{p_0^*} \text{ is a pole of } F_x(z) \\ \text{If } z_0 \text{ is a zero of } F_x(z) &\Rightarrow \frac{1}{z_0^*} \text{ is a zero of } F_x(z) \end{aligned}$$

*Proof.* Since  $F_x(z)$  is rational, there exist two finite order polynomials  $N(z)$  and  $D(z)$  such that:

$$F_x(z) = \frac{N(z)}{D(z)}$$



For this form of  $F_x(z)$  we have,

$$F_x^*(1/z^*) = \frac{N^*(\frac{1}{z^*})}{D^*(\frac{1}{z^*})}$$

and since  $F_x(z) = F_x^*(1/z^*)$  the following has to hold:

$$\frac{N(z)}{D(z)} \frac{D^*(\frac{1}{z^*})}{N^*(\frac{1}{z^*})} = 1$$

Hence there has to be perfect cancellations of all poles of the left-hand side to all its zeros. Since  $N(z)$  and  $D(z)$  are coprime, they cannot have pole-zero cancellation and since the order of  $N(z)$  is identical to the of  $N^*(\frac{1}{z^*})$ , a zero of  $N(z)$  has to be a zero of  $N^*(\frac{1}{z^*})$ . Hence,

$$\text{If } z_0 \text{ is a zero of } N(z) \Rightarrow z_0 \text{ is a zero of } N^*(\frac{1}{z^*})$$

or

$$N^*(\frac{1}{z_0^*}) = 0 \Rightarrow N(\frac{1}{z_0^*}) = 0$$

The same proof holds for the poles. □

Based on this Lemma, the Assumptions imposed on the function  $F_x(z)$  on page 128 and the stability of the spectral factor, now show that  $F_x(z)$  must have the following properties:

1. For each pole (zero) of  $F_x(z)$  there is a conjugate reciprocal pole (zero).
2. Since  $F_x(e^{j\omega}) > 0$ , it cannot be that  $F_x(z)$  has a zero on the unit circle.
3.  $F_x(z)$  does not have a pole on the unit circle. The consequence is that the spectral factor  $Q(z)$  that satisfies (7.5) also cannot have a pole on the unit circle.

This combined with the requirement for uniqueness that  $Q(\infty) = 1$  leads to a constructive solution. Before stating this solution via an illustration in Example 7.2, we state what the format of  $Q(z)$  needs to be in order to meet the requirement  $Q(\infty) = 1$ . Let the poles and zeros of  $Q(z)$  be denoted by  $p_i$  and  $z_i$  resp, then  $Q(\infty) = 1$  is met if it is written as,

$$Q(z) = \frac{\prod_{i=1}^q (1 - z_i z^{-1})}{\prod_{i=1}^p (1 - p_i z^{-1})}$$

The constructive solution is given by the following recipe to solve the Spectral factorization problem:

**Recipe to the Spectral Factorization problem:**

**Given:** A complex function  $F_x(z) \in \mathbb{C}$  for  $z \in \mathbb{C}$  that satisfies the 3 properties listed on page 131.

**Then do the following:**

**Step 1:** Compute the poles and zeros of  $F_x(z)$ . Let these be denoted resp. as  $p_i$  and  $z_i$ .

**Step 2:** Separate the stable poles and zeros from their unstable counterparts. Let these be denoted resp. as  $p_i^s$  for  $i = 1 : p$  and  $z_i^s$  for  $i = 1 : q$ . The stable ones are by the above properties  $F_x(z)$  strictly inside the unit circle. Therefore a perfect separation between the stable ones (located strictly inside the unit circle) and the unstable ones (located strictly outside the unit circle) is possible.

**Step 3:** The stable poles  $p_i^s$  and zeros  $z_i^s$  define the spectral factor as:

$$Q(z) = \frac{\prod_{i=1}^q (1 - z_i^s z^{-1})}{\prod_{i=1}^p (1 - p_i^s z^{-1})} \quad (7.7)$$

**Step 4:** Finally we find the scalar factor  $\sigma$  in the factorization given in (7.5). As we know that for  $z$  on the unit circle the quantity  $F_x(z)$  is real, we find  $\sigma$  e.g. from the following equation:

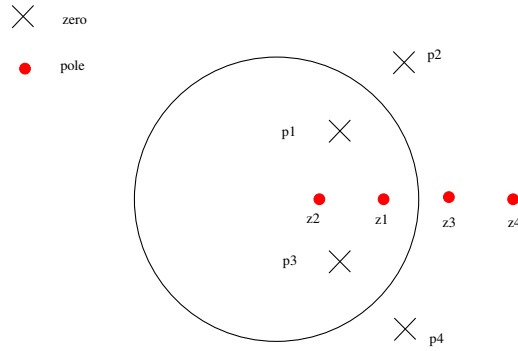
$$F_x(z_0) = \sigma Q(z_0) Q^*(1/z_0^*)$$

for  $z_0$  an arbitrary complex number on the unit circle. As an example the value  $z_0 = 1$  could be taken.

---

### Example 7.2 (Illustration Solution Spectral Factorization Problem)

As an illustration of a pole-zero pattern of function that satisfies the 3 properties on page 131, we refer to Figure 7.2. From the figure the conjugate-reciprocity



**Figure 7.2:** The poles and zeros of a complex function  $F_x(z)$  satisfying the 3 properties on page 131. The  $\times$  denotes a pole and the  $\circ$  denotes a zero.

of the poles and zeros of  $F_x(z)$  as outlined in Lemma 7.1 can clearly be observed. Following step 2 of the given recipe the stable poles are:

$$p_0, p_0^*$$

and the stable zeros are:

$$z_1, z_2$$

Therefore in step 3 we have the spectral factor given as:

$$Q(z) = \frac{(1 - z_1 z^{-1})(1 - z_2 z^{-1})}{(1 - p_0 z^{-1})(1 - p_0^* z^{-1})}$$

---

**Remark 7.2.** As stated in the Problem Formulation in subsection 7.3.1, we could start the spectral factorization from the real function  $F_x(z)$  for  $z = e^{j\omega}$ . Then the above recipe could be followed by making the substitution  $z = e^{j\omega}$ .

### 7.3.3 Use of the Spectral Factorization

When the spectral factorization has been computed the results can be used in two ways. This is explained for the spectral factorization as denoted in (7.5).

**Simulation:** When we generate a ZMWN with variance  $\sigma$ , then filtering this signal by the spectral factor  $Q(z)$  gives a filtered signal that has exactly the same Power spectrum as the given complex function  $F_x(z)$ . In this way from a given Power spectrum we can generate realizations of a stochastic process in order to do simulation studies.

**Whitening Filter:** As the spectral factor was made minimum phase, the inverse of the spectral factor is stable as well. Hence when we have a realization of the stochastic process  $x(n)$  with Power spectrum  $F_x(z)$ , the filter operation defined in Figure 7.3 delivers a signal  $v'(n)$  that is a ZMWN with variance  $\sigma$ . This whitening filter plays a crucial role in optimal filtering as will be



**Figure 7.3:** The whitening filtering given by the inverse of the spectral factor  $Q(z)$  in a spectral factorization given as in (7.5).

demonstrated in Chapter 9.

The Spectral Factorization has however many more applications than the above two illustrations. One application will be discussed in the next section and another in Chapter 9.

---

#### Example 7.3 (Whitening Filter)

Let the  $q(n)$  be a ZMWN RP with unit variance, i.e.  $\sigma_q^2 = 1$  and let  $x(n)$  be a WSS RP with Autocorrelation function given as:

$$r_x(k) = \frac{4}{3} \left( \frac{1}{2} \right)^{|k|}$$

Then these two RPs define the WSS RP  $y(n)$  as,

$$y(n) = x(n) + q(n)$$

The random processes  $x(n)$  and  $q(n)$  are assumed to be orthogonal. Now the task is to find a filter with transfer function  $H(z)$  such that when it filters  $y(n)$  its output is ZMWN.

For this purpose we first seek the z-transform of the Auto-correlation function  $r_y(k)$  of  $y(n)$ . Let this z-transform be denoted as  $P_y(z)$ , then under the assumption of orthogonality between  $x(n)$  and  $q(n)$ , we have that

$$r_y(k) = \frac{4}{3} \left( \frac{1}{2} \right)^{|k|} + \Delta(k) \quad (7.8)$$

$$\begin{aligned} P_y(z) &= \frac{1}{(1 - 0.5z^{-1})(1 - 0.5z)} + 1 \\ &= \frac{2.25 - 0.5z^{-1} - 0.5z}{(1 - 0.5z^{-1})(1 - 0.5z)} \end{aligned} \quad (7.9)$$

As a second step we compute the *spectral factorization* of  $P_y(z)$ . For that purpose we evaluate the roots of the numerator of  $P_y(z)$ . These are:

$$\frac{4.5 \pm \sqrt{4.5^2 - 4}}{2}$$

This yields (up to 4 digits) 0.2344 and  $1/0.2344$ . The spectral factorization is then given as:

$$P_y(z) = \sigma \frac{(1 - 0.2344z^{-1})(1 - 0.2344z)}{(1 - 0.5z^{-1})(1 - 0.5z)}$$

Taking  $z = 1$  yields  $\sigma = 2.1326$ . Finally, with this factorization the *whitening filter*  $H(z)$  equals:

$$H(z) = \frac{1}{\sqrt{\sigma}Q(z)} = \frac{1}{\sqrt{2.1326}} \frac{(1 - 0.5z^{-1})}{(1 - 0.2344z^{-1})}$$

## 7.4 Finding the Shaping filter of a stochastic process given its Auto-correlation Function

This Section can be considered as the time counterpart of the Spectral Factorization Problem treated in Section 7.3. As a matter of fact starting from the same motivating example given in the subsection 7.2 we can use the inverse Fourier transform on the Power Spectrum (now given as a function of the frequency variable  $\omega$ ) to obtain the Auto-correlation function.

The inverse problem based on the Auto-correlation function of a stochastic process is formulated based on the correspondence between ARMA models and Auto-correlation functions as outlined in subsection 6.4.3. As in this subsection we will treat ARMA, AR and MA model types.

### 7.4.1 Problem Formulation

The problem formulation will first be given for the most general ARMA model and then specialised to AR and MA models. In solving these inverse problems, we start however with the AR formulation, followed by the MA formulation and then closing with the solution for ARMA models.

#### The inverse ARMA problem

Assume a real-function  $r_x(k) \in \mathbb{R}$  and assume that this function is an auto-correlation function of a WSS stochastic process that has been generated by an ARMA model of the form as given in (6.10), with now the following model parameters to be unknown:

$$p, q, a(1), \dots, a(p), b(0), \dots, b(q) \quad (7.10)$$

The inverse ARMA-problem is to determine these unknown parameters from  $r_x(n)$ .

With the parameters in (7.10) estimated from the given function  $r_x(k)$ , we have a parametric model of the shaping filter that we can hopefully use as the spectral factor outlined in subsection 7.3.3. The use of the solution will be discussed for each model case (AR, MA or ARMA) separately.

**Remark 7.3.** *The variance of the ZMWN signal in the ARMA model is assumed to be 1 as this parameter can easily be integrated in the coefficients  $b(j)$ .*

#### The inverse AR problem

Assume a real-function  $r_x(k) \in \mathbb{R}$  and assume that this function is an auto-correlation function of a WSS stochastic process that has been generated by an AR model given as the difference equation,

$$x(n) + a(1)x(n-1) + \dots + a(p)x(n-p) = b(0)v(n) \quad (7.11)$$

However now again the model parameters:

$$p, a(1), \dots, a(p), b(0) \quad (7.12)$$

are not known. The inverse AR-problem is to determine these unknown parameters from  $r_x(n)$ .

#### The inverse MA problem

Assume a real-function  $r_x(k) \in \mathbb{R}$  and assume that this function is an auto-correlation function of a WSS stochastic process that has been generated by a MA model given as the difference equation,

$$x(n) = b(0)v(n) + b(1)v(n-1) + \dots + b(q)v(n-q) \quad (7.13)$$

However, again the model parameters:

$$q, b(0), \dots, b(q) \quad (7.14)$$

are not known. The inverse MA-problem is to determine these unknown parameters from  $r_x(n)$ .

## 7.4.2 Solution

In the solution we will assume the orders  $p$  and  $q$  for the different models to be known and only address the determination of the coefficients. Exercise 7.6 calls for a possible approach to detect the order  $p$  of an AR model.

The solution will make use of the Yule-Walker equation for the different models as presented in subsection 6.4.3. Taken the Remark 7.3 into account, we will consider this equation for  $\sigma_v = 1$ .

For the sake of simplicity we start treating the AR case first.

### To the inverse AR-problem

Recall the Yule-Walker equation (6.23) for AR-models, now for  $\sigma_v^2 = 1$  repeated here as,

$$r_x(n) + \sum_{\ell=1}^p a(\ell)r_x(n-\ell) = |b(0)|^2\delta(n) \quad n \geq 0 \quad (7.15)$$

Remark that for  $n > 0$ , this equation can be written only in terms of the  $a(i)$  coefficients from  $r_x(n)$ . This observation leads to a way to determine the latter coefficients. Let  $n$  in (7.15) be taken equal to  $1, 2, \dots, p$ , then we get the following set of equations:

$$\begin{bmatrix} r_x(1) & r_x(0) & r_x^*(1) & \cdots & r_x^*(p-1) \\ r_x(2) & r_x(1) & r_x(0) & \cdots & r_x^*(p-2) \\ \vdots & & & \ddots & \\ r_x(p) & r_x(p-1) & r_x(p-2) & \cdots & r_x(0) \end{bmatrix} \begin{bmatrix} 1 \\ a(1) \\ a(2) \\ \vdots \\ a(p) \end{bmatrix} = 0$$

Observing the 1 in the vector containing the parameters  $a(i)$ , we can also write this set of equations as,

$$\begin{bmatrix} r_x(0) & r_x^*(1) & \cdots & r_x^*(p-1) \\ r_x(1) & r_x(0) & \cdots & r_x^*(p-2) \\ \vdots & & \ddots & \\ r_x(p-1) & r_x(p-2) & \cdots & r_x(0) \end{bmatrix} \begin{bmatrix} a(1) \\ a(2) \\ \vdots \\ a(p) \end{bmatrix} = - \begin{bmatrix} r_x(1) \\ r_x(2) \\ \vdots \\ r_x(p) \end{bmatrix} \quad (7.16)$$

For  $r_x(n)$  given for  $n = 0, 1, \dots, p$ , if this set of equations has a unique solution, we can determine the parameters  $a(i)$  for  $i = 1, \dots, p$ . Exercise 7.5 calls for an answer when uniqueness can be guaranteed.

The full solution of the inverse AR-problem, requires to determine  $b(0)$ . Knowing the parameters  $a(i)$  for  $i = 1, \dots, p$  we consider the Yule-Walker equation (7.15) for  $n = 0$ . This results in

$$r_x(0) + \sum_{\ell=1}^p a(\ell)r_x(-\ell) = |b(0)|^2$$

or

$$r_x(0) + \sum_{\ell=1}^p a(\ell)r_x^*(\ell) = |b(0)|^2$$

As the left hand side is known, we can take  $b(0)$  as the positive root of  $r_x(0) + \sum_{\ell=1}^p a(\ell)r_x^*(\ell)$ . The value of  $b(0)$  is not required, as we only need to specify the variance of the white noise signal  $v(n)$ .

**Remark 7.4** (Use solution inverse AR-Problem). *The result of the inverse AR-Problem can be used in a similar way as the spectral factor was used, as outlined in subsection 7.3.3. Knowing the coefficients in the difference equation (7.15) and when generating a zero-mean white noise sequence with unit variance, we can simulate a realization of a stochastic process that has the given function  $r_x(n)$  as its Auto-correlation function. Further when a realization of a stochastic process with Auto-correlation would have been given the following moving average simulation model:*

$$v'(n) = \frac{1}{b(0)}x(n) + \frac{a(1)}{b(0)}x(n-1) + \cdots + \frac{a(p)}{b(0)}x(n-p)$$

*yields the whitening signal  $v'(n)$  having unit variance.*

*This results implies that the stochastic process  $x(n)$  was generated by the found AR-model.*

### To the inverse MA-problem

Recall the Yule-Walker equation (6.24) for MA-models, now for  $\sigma_v^2 = 1$  repeated as,

$$r_x(n) = \begin{cases} \sum_{\ell=n}^q b(\ell)b^*(\ell-n) & : 0 \leq n \leq q \\ 0 & : n > q \end{cases} \quad (7.17)$$

In order to see how to retrieve the  $b(j)$  coefficients from the given coefficients  $r_x(n)$  via *spectral factorization*, we consider the z-transform of (7.17). This result is given in the next lemma. Here we will make use of the *embedding* of the finite-time series  $\{r_x(n)\}_{n=-q}^q$  into a corresponding  $\infty$  time series  $r_{x,\infty}(n)$  as follows, and do likewise for the series  $\{b(j)\}_{j=0}^q$ :

$$r_{x,\infty}(n) = \begin{cases} r_x(n) & n = -q : q \\ 0 & \text{otherwise} \end{cases} \quad b_{\infty}(j) = \begin{cases} b(j) & j = 0 : q \\ 0 & \text{otherwise} \end{cases} \quad (7.18)$$

Based on this embedding we can prove the following Lemma.

**Lemma 7.5** (Spectral Factorization for the Inverse MA-problem). *Consider the Yule-Walker equations for the MA model given as (7.17). With the coefficients  $r_x(n)$  for  $n = -q : q$  and the coefficients  $b(j)$  for  $j = 0 : q$  we define the z-transforms:*

$$P_x(z) = \sum_{n=-q}^q r_x(n)z^{-n} \quad B(z) = \sum_{n=0}^q b(n)z^{-n}$$

*then the following relationship holds between these z-transforms:*

$$P_x(z) = B(z)B^*(1/z^*)$$

*Proof.* With a change of variables, we can write (7.17) also as,

$$r_x(n) = \begin{cases} \sum_{m=0}^{q-n} b(m+n)b^*(m) & : 0 \leq n \leq q \\ 0 & : n > q \end{cases}$$

Based on the embedded series  $b_\infty(j)$  and  $r_{x,\infty}(n)$ , the above expression can be written as,

$$r_{x,\infty}(n) = \sum_{m=-\infty}^{\infty} b_\infty(m+n)b_\infty^*(m)$$

Taking the z-transform and using the properties as summarized in Table 2.3, yields the following relationship between the z-transforms  $P_{x,\infty}(z)$  and  $B_\infty(z)$  of the sequences  $r_{x,\infty}(n)$  and  $b_\infty(n)$  resp,

$$\begin{aligned} P_{x,\infty}(z) &= \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} b_\infty(m+n)b_\infty^*(m)z^{-n} \\ &= \sum_{m=-\infty}^{\infty} b^*(m) \sum_{n=-\infty}^{\infty} b_\infty(m+n)z^{-n} \\ &= \sum_{m=-\infty}^{\infty} b^*(m)z^m B_\infty(z) \\ &= B_\infty^*(1/z^*)B_\infty(z) \end{aligned}$$

Since we have that

$$P_x(z) = P_{x,\infty}(z) \quad B(z) = B_\infty(z)$$

the proof is completed.  $\square$

Lemma 7.5 suggests a unique solution for an MA model that generates the given Auto-correlation data via the Spectral Factorization. This is summarized in the following recipe.

**Recipe to the inverse MA-problem:**

**Given:** The index  $q$  and the auto-correlation function  $r_x(n) \in \mathbb{C}$  for  $n = -q : q$ .

**Then do the following:**

**Step 1:** Determine the z-transform  $P_x(z)$  of the given sequence  $r_x(n)$  as in Lemma 7.5.

**Step 2:** Perform a spectral factorization of  $P_x(z)$  denoted as:

$$P_x(z) = \alpha Q(z)Q^*(1/z^*)$$

**Step 3:** The transfer function  $H(z)$  of an MA model is given as:

$$H(z) = \sqrt{\alpha}Q(z) \tag{7.19}$$

**Remark 7.6** (Use solution inverse AR-Problem). *As in Remark 7.4, we can use the result of the inverse MA-Problem e.g. to generate an artificial sequence with the given Auto-correlation function. Or we can generate a whitening sequence given a realization*



of the stochastic process  $x(n)$ . For that purpose we write the transfer function  $H(z)$  in (7.19) explicitly as:

$$H(z) = \beta(0) + \beta(1)z^{-1} + \cdots + \beta(q)z^{-q}$$

then filtering  $x(n)$  with the AR-model with transfer function:

$$\frac{1/\beta(0)}{1 + \frac{\beta(1)}{\beta(0)}z^{-1} + \cdots + \frac{\beta(q)}{\beta(0)}z^{-q}}$$

delivers a ZMWN with unit variance. It could be checked that the poles of this AR-model are within the unit circle. Why?

### To the inverse ARMA-problem

The solution to the inverse ARMA-problem will be given in **two** steps:

1. The computation of the coefficients  $a(i)$  for  $i = 1 : p$ .
2. Using the computed coefficients  $a(i)$  for  $i = 1 : p$ , then the  $b(j)$  for  $j = 0, q$  are determined.

Again we make the assumption that  $\sigma_v^2 = 1$ . These two computational steps are subsequently described in a way that easily leads to a recipe to solve the inverse ARMA-problem. For the first step this is the solution to a set of equations. The second step is more involved and again leads to a Spectral Factorization. As it is more involved the second step is subdivided in three sub-steps.

**Step 1 — Computing the coefficients  $a(i)$  for  $i = 1 : p$ :** Recall the Yule-Walker equation (6.21) for ARMA-models, now for  $\sigma_v^2 = 1$  repeated as,

$$r_x(n) + \sum_{\ell=1}^p a(\ell)r_x(n-\ell) = \begin{cases} \sum_{\ell=n}^q b(\ell)h^*(\ell-n) & : 0 \leq n \leq q \\ 0 & : n > q \end{cases} \quad (7.20)$$

Then inspired by the solution to the inverse AR-problem, we take  $n = q+1, \dots, q+p$  to obtain the following set of equations in the unknown coefficients  $a(i)$  for  $i = 1 : p$ .

$$\begin{bmatrix} r_x(q) & \cdots & r_x(q-p+1) \\ \vdots & \ddots & \\ r_x(q+p-1) & \cdots & r_x(q) \end{bmatrix} \begin{bmatrix} a(1) \\ \vdots \\ a(p) \end{bmatrix} = - \begin{bmatrix} r_x(q+1) \\ \vdots \\ r_x(q+p) \end{bmatrix} \quad (7.21)$$

These are called the *modified Yule-Walker equation*.

**Step 2 — Computing the coefficients  $b(j)$  for  $j = 0 : q$ :** Knowing the coefficients  $a(i)$  we can compute from the Yule-Walker equation (7.20) for  $n = 0 : q$ , and determine the coefficients  $c(j)$  for  $j = 0 : q$  as follows:

$$\begin{bmatrix} r_x(0) & r_x^*(1) & \cdots & r_x^*(p) \\ r_x(1) & r_x(0) & \cdots & r_x^*(p-1) \\ \vdots & & \ddots & \\ r_x(q) & r_x(q+1) & \cdots & r_x(0) \end{bmatrix} \begin{bmatrix} 1 \\ a(1) \\ \vdots \\ a(p) \end{bmatrix} = \begin{bmatrix} c(0) \\ c(1) \\ \vdots \\ c(q) \end{bmatrix} \quad (7.22)$$

Via the Yule-Walker equation (7.20) the relation between these now known coefficients and the unknown coefficients  $b(j)$  and the unknown impulse response coefficients  $h(j)$  is given as follows:

$$c(k) = \sum_{\ell=k}^q b(\ell)h^*(\ell - k) \quad (7.23)$$

Inspired by the approach to solve the inverse MA-problem, we consider the z-transform of (7.23). For this reason we do the following embedding:

$$c_{\infty}(n) = \begin{cases} c(n) & n = 0 : q \\ 0 & n > q \\ \text{unknown} & n < 0 \end{cases} \quad b_{\infty}(j) = \begin{cases} b(j) & j = 0 : q \\ 0 & \text{otherwise} \end{cases}$$

$$h_{\infty}(j) = \begin{cases} h(j) & j \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

As in Lemma 7.5, (7.23) can be used to derive the following relation between the corresponding z-transforms of the above  $\infty$  length series:

$$C_{\infty}(z) = B_{\infty}(z)H_{\infty}^*(1/z^*) \quad (7.24)$$

Due to the unknown values of the  $c_{\infty}(n)$  coefficients for  $n < 0$ , only the causal part of the  $C_{\infty}(z)$  is known. This causal part is with the notation introduced on page 25 denoted as  $[C_{\infty}(z)]_+$ . So the sub-problem to be addressed now is to determine from  $[C_{\infty}(z)]_+$  the coefficients of the numerator polynomial of an ARMA model. This is done in the following Lemma and Theorem.

**Lemma 7.7** (A Power Spectrum is defined by its Causal part only). *Let a complex function (like a Power Spectrum)  $P_x(z)$  satisfy,*

$$P_x(z) = P_x^*(1/z^*)$$

*and let this function have a z-transform denoted as:*

$$P_x(z) = \sum_{n=-\infty}^{\infty} p_n z^{-n}$$

*then this function is fully defined if  $[P_x(z)]_+$  is given.*

*Proof.* Exercise 7.8 calls for a proof of this lemma. □

**Theorem 7.8** (Solution second step inverse-ARMA problem). *Let the following conditions hold,*

1. *the unknown ARMA model in the inverse ARMA-problem have a causal transfer function denoted as:*

$$H_{\infty}(z) = \frac{B(z)}{A(z)}$$

*with  $B(z)$  a polynomial in  $z \in \mathbb{C}$  of order  $q$  and  $A(z)$  a polynomial of order  $p$ ,*

2. the denominator polynomial  $A(z)$  be given,
3. let the relation in (7.24) be given with the quantities in that equation as defined above the relation.
4. Let the complex function  $P_y(z)$  satisfy the conditions of Lemma 7.7 be defined by its causal part given as:

$$[P_y(z)]_+ = [[C_\infty(z)]_+ A^*(1/z^*)]_+$$

5. Let  $P_y(z)$  have a spectral factorization denoted as:

$$P_y(z) = \alpha Q(z) Q^*(1/z^*)$$

then,

1.  $P_y(z)$  is a polynomial in  $z$ , and,
2. the numerator polynomial  $B(z)$  of the unknown ARMA model is given as:

$$B(z) = \sqrt{\alpha} Q(z)$$

*Proof.* Using the expression of  $H_\infty(z)$  in (7.24) and using the fact that  $B_\infty(z) = B(z)$  (see Proof of Lemma 7.5), (7.24) becomes:

$$C_\infty(z) = B(z) \frac{B^*(1/z^*)}{A^*(1/z^*)}$$

Multiplying both sides with  $A^*(1/z^*)$  yields:

$$C_\infty(z) A^*(1/z^*) = B(z) B^*(1/z^*) \quad (7.25)$$

If we denote the above complex function in  $z$  as  $P_y(z)$ , then it can be concluded that this function satisfies the conditions of Lemma 7.7 and it is polynomial. Let  $C_\infty(z)$  be decomposed into its causal and anti-causal part as:

$$C_\infty(z) = [C_\infty(z)]_+ + [C_\infty(z)]_-$$

(with its causal part known).

Then the causal part of the function  $P_y(z)$  is given by the expression:

$$[P_y(z)]_+ = \left[ [C_\infty(z)]_+ A^*(1/z^*) + \underbrace{[C_\infty(z)]_- A^*(1/z^*)}_{\text{anti-causal}} \right]_+$$

The under-braced part in this expression is purely anti-causal. Why is that? Therefore,

$$[P_y(z)]_+ = [[C_\infty(z)]_+ A^*(1/z^*)]_+$$

As  $P_y(z)$  satisfies Lemma 7.7 this part fully defines  $P_y(z)$  and with its spectral factorization given as in the Theorem, we find from (7.25) that the numerator  $B(z)$  can be taken as indicated in the Theorem.  $\square$

---

**Example 7.4 (Inverse ARMA problem)**

In this example we are given the following three values of the Auto-correlation function of WSS RP  $x(n)$ :

$$r_x(0) = 26 \quad r_x(1) = 7 \quad r_x(2) = 3.5$$

As a first step to determine the ARMA model that would generate a WSS RP by filtering ZMWN with unit variance and that has the same Auto-correlation function, we determine the coefficients of the denominator of this ARMA model via the Yule-Walker equations (7.21) for  $p = 1$ . Then we find that  $a(1) = -\frac{1}{2}$ , and as such,

$$A(z) = 1 - \frac{1}{2}z^{-1}$$

For the second step, we compute the coefficients  $b(0)$  and  $b(1)$ . We start this calculation by computing the polynomial  $[C(z)]_+$ . Using the computed coefficient  $a(1)$  and the Yule-Walker equations (7.22) (again for  $p = 1$ , we find,

$$[C(z)]_+ = 22.5 - 6z^{-1}$$

This defines the polynomial  $P_y(z)_+$  as in condition 4. of Theorem 7.8 as,

$$\begin{aligned} [P_y(z)]_+ &= \left[ [C(z)]_+ A^*(1/z^*) \right]_+ \\ &= \left[ \left( \frac{45}{2} - 6z^{-1} \right) \left( 1 - \frac{1}{2}z \right) \right]_+ \\ &= \left[ -\frac{45}{4}z + \left( \frac{45}{2} + \frac{6}{2} \right) - 6z^{-1} \right]_+ \\ &= \frac{51}{2} - 6z^{-1} \end{aligned}$$

Then following Lemma 7.7,  $P_y(z)$  now becomes  $P_y(z) = -6z + \frac{51}{2} - 6z^{-1}$ . The coefficients  $b(0)$  and  $b(1)$  can now be found via a *spectral factorization* of  $P_y(z)$ . This spectral factorization is given as,

$$P_y(z) = 24 \left( 1 - \frac{1}{4}z^{-1} \right) \left( 1 - \frac{1}{4}z \right)$$

From this it follows that,

$$B(z) = 2\sqrt{6} \left( 1 - \frac{1}{4}z^{-1} \right)$$

And the ARMA model that we seek is given by,

$$\frac{1 - \frac{1}{2}z^{-1}}{2\sqrt{6} \left( 1 - \frac{1}{4}z^{-1} \right)}$$

---

---

## References

- [1] N.J. Usdin, *Discrete Simulation of Colored Noise and Stochastic Processes and  $1/f^\alpha$  Power Law Noise Generation*. Proc. IEEE, Vol. 83:5, 1995.
- [2] V.A. Ambartsumian, *On the Relationship between the Solution and the Resolvente of the Integral Equation of the Radiative Balance*. Zeitschrift für Physik, Vol. 52:3-4, pp. 263-267, 1929 (In German).
- [3] G. Volpe and G. Volpe, *Simulation of a Brownian particle in an optical trap*, American J. Phys, Vol. 81(3), pp. 224–230, 2013.
- [4] L. Ljung, *System Identification: Theory for the User*. Upper Saddle River, New Jersey: Prentice Hall, second ed., 1999.
- [5] M. Verhaegen and V. Verdult, *Filtering and System Identification: A Least Squares Approach*. Cambridge: Cambridge University Press, 2007.
- [6] M. Verhaegen, G. Vdovin and O. Soloviev, *Control for High Resolution Imaging*. Lecture Notes, Delft University of Technology (Course sc4045), 2015.
- [7] D.A. Duev, et. al., *Spacecraft VLBI and Doppler tracking: algorithms and implementation*, Astronomy & Astrophysics manuscript No. AA18885-12, 2012.
- [8] D. Clay, *Linear Algebra and its Applications*, 4th-Ed., Addison-Wesley, 2012.

---

## Exercises

**Exercise 7.1** Consider the second rational function  $F_x(z)$  considered in Example 7.1, repeated here:

$$F_x(z) = \frac{az + b + cz^{-1}}{dz + e + fz^{-1}} \left( = \frac{a + bz^{-1} + cz^{-2}}{d + ez^{-1} + fz^{-2}} \right)$$

with  $a, b, c, d, e, f \in \mathbb{R}$ . When  $F_x(z)$  satisfies the symmetry property,

$$F_x(z) = F_x^*(1/z^*)$$

and the coefficients satisfy,

$$a \neq c \quad d \neq f \quad \text{and} \quad af \neq 0$$

then show that the following relationship holds between the coefficients of  $F_x(z)$ :

$$\frac{a}{d} = \frac{c}{f} = \frac{b}{e}$$

**Exercise 7.2** Determine the spectral factorization of the complex function:

$$F_x(z) = \frac{1 - 2.5z^{-1} + z^{-2}}{1 - 2.05z^{-1} + z^{-2}}$$

**Exercise 7.3** For the given real function  $F_x(\omega)$ :

$$F_x(\omega) = \frac{25 - 24\cos\omega}{26 - 10\cos\omega}$$

for  $\omega \in [0, \pi]$ .

- (a) Determine the spectral factor  $Q(z)$  and the scalar  $\sigma$  in the spectral factorization in (7.5).
- (b) Determine the whitening filter  $H(z)$  such that when filtering a given realization of stochastic process that has the Power spectrum  $F_x(\omega)$  by  $H(z)$  we obtain ZMWN with unit variance.

**Exercise 7.4** Consider the (real) Auto-correlation function given by:

$$r_x(k) = \frac{4}{3} \left(\frac{1}{2}\right)^{|k|}$$

then determine the parameter  $a$  and  $b$  in the AR model:

$$x(n) + ax(n-1) = bv(n) \quad v(n) \text{ is ZMWN}(1)$$

**Exercise 7.5** What is the condition on the matrix,

$$\begin{bmatrix} r_x(0) & r_x^*(1) & \cdots & r_x^*(p-1) \\ r_x(1) & r_x(0) & \cdots & r_x^*(p-2) \\ \vdots & & \ddots & \\ r_x(p-1) & r_x(p-2) & \cdots & r_x(0) \end{bmatrix}$$

in the set of equations (7.16) to determine the  $a(i)$  for  $i = 1, \dots, p$  in the inverse AR-problem has a *unique* solution.

**Exercise 7.6** Consider the solution to the inverse AR-problem. Then show,

- (a) for the general case assuming order  $p$ , the following expression holds,

$$\begin{bmatrix} r_x(1) & r_x(0) & r_x^*(1) & \cdots & r_x^*(p-1) \\ r_x(2) & r_x(1) & r_x(0) & \cdots & r_x^*(p-2) \\ \vdots & & & \ddots & \\ r_x(p) & r_x(p-1) & r_x(p-2) & \cdots & r_x(0) \\ r_x(p+1) & r_x(p) & r_x(p-1) & & r_x(1) \\ r_x(p+2) & r_x(p+1) & r_x(p) & & r_x(2) \\ \vdots & & \ddots & & \vdots \\ r_x(p+\alpha) & r_x(p+\alpha-1) & r_x(p+\alpha-2) & \cdots & r_x(\alpha) \end{bmatrix} \begin{bmatrix} 1 \\ a(1) \\ a(2) \\ \vdots \\ a(p) \end{bmatrix} = 0$$

for  $\alpha > 1$ .

- (b) if the “unknown” AR-model that has generated the Auto-correlation function  $r_x(n)$  has now order  $p = 2$ , then show that the matrix,

$$\begin{bmatrix} r_x(2) & r_x(1) & r_x(0) \\ r_x(3) & r_x(2) & r_x(1) \\ r_x(4) & r_x(3) & r_x(2) \\ r_x(5) & r_x(4) & r_x(3) \end{bmatrix}$$

has rank 2. [For a definition of the rank of a matrix we refer to the linear algebra book of [8].]

**Exercise 7.7** Consider the MA model with coefficients:

$$b(0) = 1 \quad b(1) = -2 \quad q = 1$$

- (a) Determine the transfer function of the MA model and denote this by  $B(z)$
- (b) What is the zero of  $B(z)$ ? Is it stable?
- (c) Determine the Auto-correlation function of the output  $x(n)$  of this MA model for the variance of its input equal to 1. Denote this Auto-correlation function as  $r_x(n)$ .
- (d) Given  $r_x(n)$ , determine the solution of the inverse MA-problem. Denote the transfer function of this MA model by  $H(z)$ .
- (e) What is the zero of  $H(z)$ ?

**Exercise 7.8** Prove Lemma 7.7.

**Exercise 7.9** Let the following coefficients of the Auto-correlation function of an ARMA model with  $p = 1$  and  $q = 1$  be given:

$$r_x(0) = 8.4 \quad r_x(1) = 1.6 \quad r_x(2) = 1\frac{1}{15}$$

then determine the parameters of an ARMA model with transfer function:

$$\frac{b(0) + b(1)z^{-1}}{1 + a(1)z^{-1}}$$

such that the Auto-correlation function of the stochastic process that results by filtering unit white noise with this filter for lags 0,  $\pm 1$  and  $\pm 2$  corresponds to the given Auto-correlation coefficients.

**Exercise 7.10** Consider a linear, time-invariant, dynamic minimum-phase system  $G$  with the discrete-time stochastic process  $u$  as its input and the discrete-time stochastic process  $y$  as its output. Suppose that the input  $u$  is equal to a white noise process with variance 1 and that the output is given by:

$$P_y(\omega) = \frac{5 + 4 \cos(\omega)}{5 + 3 \cos(\omega)}.$$

Determine the transfer function  $G(z)$  of the system  $G$ .

**Exercise 7.11** Consider the discrete-time, wide-sense stationary (WSS) stochastic process  $x(n)$  with mean 0 and power spectral density function

$$P_x(e^{j\omega}) = 1.$$

Assume  $x(n)$  is the input of a linear, time-invariant, discrete-time system with unit sample response  $h(n)$ . Then, the output of the system is a discrete-time stochastic process  $y(n)$  that can be written as  $y(n) = h(n) \star x(n)$ , with  $\star$  the convolution operator. Furthermore, it is given that:

- $h(0) > 0$ ,
  - $h(n) = 0$  for  $n < 0$  and  $n \geq 2$ ,
  - $P_y(e^{j\omega}) = \frac{5}{4} - \cos(\omega)$ .
- (a) Give the pole-zero plot of  $P_y(z)$  and give a unit sample response  $h(n)$  that is consistent with the information given.
- (b) Show that the answer of (a) is not unique by giving a second unit sample response  $h(n)$  that is consistent with the information given.
- (c) Next, consider the process  $w(n) = g(n) \star x(n)$  with  $g(n) = 0$  for  $n < 0$  and  $n \geq 2$ , and  $g(0) > 0$ . Assume that:

$$P_w(e^{j\omega}) = \frac{5}{4} + \cos(\omega)$$

and

$$P_{yw}(e^{j\omega}) = \frac{1}{4}e^{i\omega} - e^{-i\omega}.$$

Give the pole-zero plots of  $P_w(z)$  and  $P_{yw}(z)$  and give a possible unit sample response  $h(n)$  that is consistent with the information given.

- (d) Is the answer given in (c) unique? If yes, explain. If not, give another unit sample response  $h(n)$  that is consistent with the information given.

**Exercise 7.12** Consider the linear, causal, stable, discrete-time dynamic system characterized by the following transfer function:

$$H(z) = \frac{1 - \frac{1}{2}z^{-1}}{1 - \frac{1}{3}z^{-1}}$$

Determine the unit sample response  $h(k)$  of filter  $H(z)$ .



## Chapter 8

# Parameter Estimation: The Linear Least Squares Method

**After studying this chapter you can:**

- formulate the estimation of the unknown parameters of an AR model as a Linear Least Squares (LLSQ) problem when a finite length sequence of a single realization of the output of the AR model is given
- analyse the unbiasedness and the covariance matrix of the estimate of the LLSQ solution to the estimation of the AR model parameters
- generalize the LLSQ estimation of the AR model parameters to other estimation problems and,
- use the so-called orthogonality condition to derive a solution to generic LLSQ problems.

---

## 8.1 Introduction

In Chapter 7 a number of inverse problems were considered in the time and frequency domain. In Chapter 7 the data set (observations) that was the starting point for the Power Spectrum inverse problem, treated in Section 7.3, was a rational Power Spectrum. For the Auto-correlation inverse problem, treated in Section 7.4, it was a function that could be considered as the Auto-correlation function of a particular AR, MA or ARMA model. In general such starting data is ideal and generally in practical real-life experiments collecting real-life observations, such information is generally not available. This is because the information retrieved from reality is much more complex. If we consider e.g. the example about modeling turbulence as given in Section 7.2 and Kolmogorov's  $\frac{5}{3}$  power law, as illustrated by (7.1), then it is clear that turning such a Power Spectrum into a rational function of the variable  $z = e^{jf}$  is not possible. Another source of problems is that the data collected in real-life experiments is contaminated by "noise" due to e.g. the finite accuracy of the measurement equipment. Though it may in many applications be assumed that such noise is statistically uncorrelated from the stochastic process that is under study, it complicates the solution to (realistic) inverse problems.

As indicated in Figure 7.1, Kolmogorov's  $\frac{5}{3}$  power law is an approximation of the real-life measurements. Therefore one could ask the question whether it is not better to address the inverse problems starting from real-life measurements. Referring to Figure 7.1, this would then be starting from the blue "raw" data.

Although this is a highly relevant question, the analysis and solution of such inverse problems using real-life measurements is in general much more complicated. This is why in these introductory course notes, we restrict to the analysis of a simple inverse problem to estimating the parameters of an AR model from one realization of a stochastic process that is addressed to be generated by the AR model.

This inverse problem can be formulated as a *Linear Least Squares (LLSQ)* problem, a problem that has been pioneered by C.F. Gauss, as outlined in Section 1.3. The formulation of this LLSQ problem and an illustration of it is given in Section 8.2. The solution of the problem is given in Section 8.3. Here we present this solution as the outcome of a deterministic optimization problem. As this optimization problem is quadratic in the unknown parameters, it can be solved using the so-called *orthogonality principle*. Since the LLSQ method can be and has been applied in a wide number of fields and applications, other examples are given of the application of the LLSQ method. A final topic analysed in this Chapter is the effect of dealing with finite length data sequences. This is often a "natural" limitation when working with real-life experiments, and it causes the fact that the computed parameter values are estimates. As such these estimates can be considered as random variables. For that reason we analyse via statistical notions of **bias** and the *covariance matrix* how we could quantify the quality of an estimate. This is done in Section 8.4.

---

## 8.2 The Linear Least Squares (LLSQ) Problem

### 8.2.1 The LLSQ problem for estimating the parameters of an AR model

First we provide a linguistic formulation of the problem and then rephrase it in a mathematical optimization problem.

*Consider the AR model given by (7.11) and assume that a finite number of  $N$  data points of a single realization of the stochastic process  $x(n)$  are available. Let this set be denoted as,*

$$\{x(n)\}_{n=0}^{N-1} \quad (8.1)$$

*Assume further that the order  $p$  in (7.11) is known, then the problem that we would like to analyse is to estimate the parameters  $a(i)$  for  $i = 1 : p$  in an optimal manner from the given data set (8.1).*

For the mathematical formulation we need to define the notion ‘optimality’. Here we find inspiration in the work of Gauss. Let the parameters that we seek be denoted as  $a_N(i)$  for  $i = 1 : p$  and stored in the vector  $a_N$  as:

$$a_N = \begin{bmatrix} a_N(1) \\ \vdots \\ a_N(p) \end{bmatrix}$$

Subsequently we make a “copy” of the AR part with these parameters and call that AR copy an error signal  $e(n; a_N)$ , then we arrive at the following expression:

$$x(n) + a_N(1)x(n-1) + \cdots + a_N(p)x(n-p) = e(n; a_N) \quad (8.2)$$

Here the notation  $e(n; a_N)$  emphasizes the explicit dependency of the error signal  $e(n)$  on the parameter values  $a(j)$  for  $j = 1 : p$ .

The least squares approach now aims at finding estimates for the parameters  $a_N(i)$  by minimizing an ergodic approximation based on finite data points of the variance of the error signal  $e(n; a_N)$ . It should be remarked that when  $x(n)$  is a WSS process and for constant parameters  $a_N(i)$  Theorem 6.1 shows that the signal  $e(n; a_N)$  is also WSS and therefore the assumption that it is (Auto-correlation) ergodic is appropriate.

That minimizing the variance of the error signal  $e(n; a_N)$  allows to find the true AR parameters is highlighted in the following Example.

---

#### Example 8.1 (Variance Minimization)

Consider the AR model for  $p = 1$  given as:

$$x(n) - ax(n-1) = v(n) \quad |a| < 1$$

and  $v(n)$  ZMWN with variance equal to 1. With the stochastic process generated in this way, we define as in (8.2) for general  $p$  a new stochastic process  $e(n)$  as follows:

$$e(n) = x(n) - (a + \epsilon)x(n-1)$$

Now we evaluate the variance of the  $e(n)$  as a function of this (offset) parameter  $\epsilon$ . Combining the above two equations allows to express  $e(n)$  as:

$$e(n) = v(n) - \epsilon x(n-1)$$

Then as asked for in Exercise 8.1 it can be shown that:

$$E[x(n-1)v^*(n)] = 0$$

As a consequence, in the same line of Exercise 3.10, the variance of the stochastic process  $e(n)$  now equals:

$$E[|e(n)|^2] = E[|v(n)|^2] + \epsilon^2 E[|x(n-1)|^2] = 1 + \epsilon^2 r_x(0)$$

As  $r_x(0)$  is non-negative, see Property 5.9, we have that for the case  $r_x(0) > 0$  the minimum value of the variance of the stochastic process  $e(n)$  is reached for

$$\boxed{\epsilon = 0} \quad (8.3)$$

The additional requirement  $r_x(0) > 0$  easily holds for the case  $x(n)$  is generated. (Please check this yourself).

---

The example 8.1 can be generalized for  $p > 1$  (see next Chapter 9) and therefore the original parameters of the AR model can be found by *minimizing the variance of the error signal*  $e(n; a_N)$  defined in (8.2).

However as we assumed to only have a finite set of data available, we have to work with a finite sample approximation of the variance of  $e(n)$ , and for that purpose we assume it to be (Auto-correlation) ergodic. By this assumption, the variance of  $e(n; a_N)$  can be approximated by its sample average, given as:

$$\hat{r}_e(0, N-p) = \frac{1}{N} \sum_{n=p}^{N-1} |e(n; a_N)|^2 \quad (8.4)$$

Based on this ergodic approximation of the variance, the *optimal* estimation of the AR parameters is now expressed by the following optimization problem,

$$\min_{a_N(i)} J(a_N(i)) = \min_{a_N(i)} \frac{1}{N} \sum_{n=p}^{N-1} |e(n; a_N)|^2$$

this can be reformulated with the definition of the error signal  $e(n; a_N)$  in (8.2), as

$$\min_{a_N(i)} J(a_N(i)) = \min_{a_N(i)} \frac{1}{N} \sum_{n=p}^{N-1} |x(n) + a_N(1)x(n-1) + \dots + a_N(p)x(n-p)|^2 \quad (8.5)$$

The solution of this optimization problem are called estimates. This is because they are random variables. The properties of these estimates will be analysed in Section 8.4.

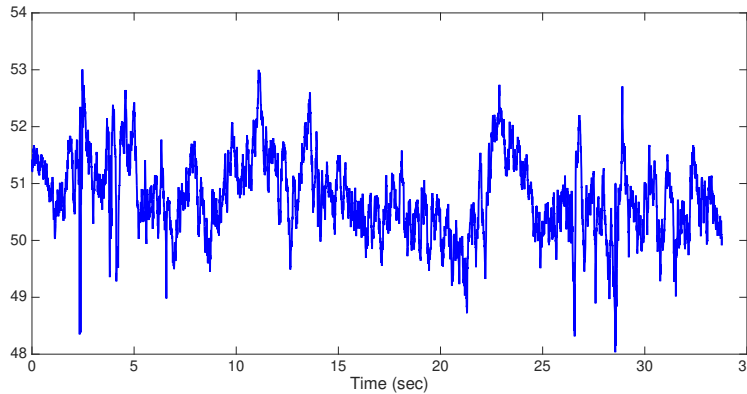
We conclude this section with a motivating example on the relevance of estimating AR models (or their more general counterparts like ARMA or state space model) from experimental data.

---

**Example 8.2 (Turbulence induced changes of the refraction index)**

Turbulence induces wavefront aberrations that reduce the image resolution of images recorded with e.g. ground based telescopes. The field of Adaptive Optics, using wavefront sensors and deformable mirrors, aims at real-time correcting these wavefront aberrations. See e.g. [3] for an introduction to this field of optics.

The aberrations of the spatial slopes of the wavefront aberrations can be measured with a Shack-Hartmann wavefront sensor. This is a two dimensional array of small lenses (see [2] for a brief technical description of this sensor). For example, on the William-Herschel telescope [1] this is an  $8 \times 8$  orthogonal micro-lens array, which partitions the telescope aperture in sub-apertures with an effective size of 0.5m. Each sub-aperture has a field of view of 2.2 arcsec and is imaged on a  $8 \times 8$  sub-region of a  $64 \times 64$  CCD camera, giving rise to a sensitivity of approximately 0.27 arcsec/pixel. At each sample instance the pixels measured over each  $8 \times 8$  sub-region are averaged via the so-called *centroiding* method [1]. This determines the spatial average of the measured intensities, called the centroid, at each sample instance. An example of a 30 sec time record of one of these centroids recorded in a real-life experiment with the telescope with a sample frequency of 296 Hz is plotted in Figure 8.1. In order to test new Adaptive Optics



**Figure 8.1:** One of the 64 centroids recorded with the Shack-Hartmann wavefront sensor during one of the experiments within the Joint Observatories Seeing Evaluation (JOSE) project with the William Herschel telescope.

control methods, prior to testing on a real telescope (which may be extremely expensive) simulation in the lab or on a computer is crucial. In order to conduct realistic simulations it may be helpful to generate from a real-life recording such as displayed in Figure 8.1 artificial recordings of arbitrary duration. One way of

doing this is to try to use the least squares method for estimating the parameters in an AR model and then using these estimated parameters together with a synthetic white noise sequence (of desirable length) to generate a realization of the stochastic process, see Remark 7.4. Inspecting the AR model in (7.11), we see that prior to estimation of the AR parameters, the order  $p$  has to be selected. This is a topic of research for which we refer to more advanced courses, such as outlined in [4]. In that course you will learn also to estimate the parameters of other model type, such as state space models.

---

## 8.2.2 Linear Regression

The LLSQ problem (8.5) for estimating the parameters can be easily made more abstract in order to demonstrate its generality. This is done by putting the problem (8.5) in a matrix-vector context. Hereby we define the following matrix and vectors:

$$\begin{aligned} X_N &= \begin{bmatrix} x(p-1) & x(p-2) & \cdots & x(0) \\ x(p) & x(p-1) & & x(1) \\ \vdots & & \ddots & \vdots \\ x(N-2) & x(N-3) & \cdots & x(N-p-1) \end{bmatrix} & y_N &= \begin{bmatrix} x(p) \\ x(p+1) \\ \vdots \\ x(N-1) \end{bmatrix} \\ a_N &= \begin{bmatrix} a_N(1) \\ \vdots \\ a_N(p) \end{bmatrix} & e_N &= \begin{bmatrix} e(p; a_N) \\ e(p+1; a_N) \\ \vdots \\ e(N-1; a_N) \end{bmatrix} \end{aligned} \quad (8.6)$$

If we in addition define the norm (or length) of a complex vector  $e_N$ , generalizing the definition for real vectors in [5], as follows:

$$\|e_N\|_2 = \sqrt{|e(p; a_N)|^2 + |e(p+1; a_N)|^2 + \cdots + |e(N-1; a_N)|^2} = \sqrt{\sum_{n=p}^{N-1} |e(n; a_N)|^2}$$

then we can write the LLSQ cost function (8.5) as,

$$J(a_N) = \frac{1}{N} \|e_N\|^2 = \frac{1}{N} \|X_N a_N + y_N\|_2^2$$

and the LLSQ problem can be formulated as,

$$\boxed{\min_{a_N} J(a_N) \quad \text{or} \quad \min_{a_N} \frac{1}{N} \|X_N a_N + y_N\|_2^2} \quad (8.7)$$

If we abstract the matrices in this LLSQ problem by defining the matrix  $X_N$  and vectors  $a_N, y_N$  differently (according to the needs of the problem at hand) the problem (8.7) is widely applicable. The problem is known under various names, such as linear regression, and belongs to the class of *quadratic optimization problems*. [6]

With the definition of the signal  $e(n; a_N)$  in the vector  $e_N$  we arrive at the following equation, which will be referred to in this book as the *data equation*

$$X_N a_N + y_N = e_N \quad (8.8)$$

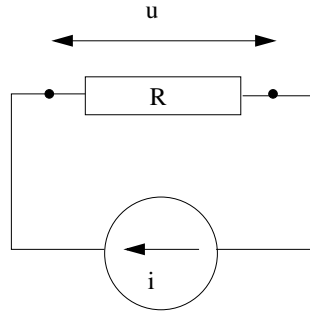
This equation contains  $N$  unknowns: the  $p$  coefficients in  $a_N$  and the  $N - p$  coefficients in  $e_N$ . Therefore by the LLSQ problem formulation we turn this into a problem in  $p$  parameters only!

We conclude this section by an example illustrating the abstraction from an AR parameter estimation problem to other estimation problems.

---

### Example 8.3 (Resistor Value Estimation)

Consider Ohms famous law  $u = Ri$ , along with the quantities defined in Figure 8.2 If we want to find the value of the resistor experimentally taking into



**Figure 8.2:** Schematic of an electrical circuit with  $u$  the voltage over the resistor with value  $R$  and current  $i$ .

consideration limitations of the measurement devices, a LLSQ problem is formulated. For that purpose we select different currents  $i(n)$  with the current source, and measure the corresponding voltage  $u(n)$  with a voltage meter. If that meter suffers from an offset  $u_0$  and random measurement errors  $e(n)$ , the measured voltage  $u_m(n)$  is assumed to be related to the true voltage as:

$$u_m(n) = u(n) + u_0 + e(n)$$

Based on Ohms law we can relate the measurements in the data equation format as:

$$u_m(n) = Ri(n) + u_0 + e(n)$$

and seek to *estimate* the resistor value  $R$  and the offset  $u_0$  via a least squares problem. Assume that we have  $N$  measurements of the voltage  $u_m(n)$  and the current  $i(n)$  for  $n = 0 : N - 1$  and the values of the resistor, offset and the error signal are resp. denoted as:

$$R_N, u_{0,N} \quad \text{and} \quad e(n; \begin{bmatrix} R_N \\ u_{0,N} \end{bmatrix})$$

Then following Ohms law, the given measurements and the above defined quantities are related as:

$$e(n; \begin{bmatrix} R_N \\ u_{0,N} \end{bmatrix}) = u_m(n) - R_N i(n) - u_{0,N}$$

and a least squares cost function for estimating the parameters  $R_N, u_{0,N}$  (assuming all quantities to be real) reads,

$$J\left(\begin{bmatrix} R_N \\ u_{0,N} \end{bmatrix}\right) = \frac{1}{N} \sum_{n=0}^{N-1} \left( e(n; \begin{bmatrix} R_N \\ u_{0,N} \end{bmatrix}) \right)^2$$

If we now store the given measurements and defined quantities in the resp. the matrix  $X_N$  and vectors  $y_N, a_N$  as follows:

$$X_N = - \begin{bmatrix} i(0) & 1 \\ i(1) & 1 \\ \vdots & \vdots \\ i(N-1) & 1 \end{bmatrix} \quad y_N = \begin{bmatrix} u_m(0) \\ u_m(1) \\ \vdots \\ u_m(N-1) \end{bmatrix} \quad a_N = \begin{bmatrix} R_N \\ u_{0,N} \end{bmatrix}$$

the estimation of the physical values of the resistor and the offset can be done via the solution of the LLSQ problems (8.7).

---

It is seen that for the AR parameter estimation problem formulation in (8.5), the problem defined in Example 8.3 that the columns of the matrix  $X_N$  and the vector  $y_N$  are defined from *sampled* signals. This is often the case when abstracting these matrix-vector quantities. In the general case the columns of the matrix  $X_N$  are referred to as the *independent variables*, while the vector  $y_N$  is called the *dependent variable*. This is because the matrix vector product  $X_N a_N$  is a linear combination of the columns of  $X_N$ . In the solution of the LLSQ problem it will become clear that its solution gives an approximation of the vector  $y_N$  that is exactly a linear combination of the columns of  $X_N$ , i.e. (linearly) dependent on  $X_N$ .

---

## 8.3 Solution to the AR parameter estimation problem

### 8.3.1 Deriving the solution

As illustrated by the solution of quadratic cost function of complex parameters, as outlined in Section 2.4 for two parameters, the necessary and sufficient conditions that characterizes the minimum of the cost function  $J(a_N)$  in (8.5) is given by:

$$\frac{\partial J(a_N)}{\partial a_N^*(k)} = 0 \quad \text{for } k = 1 : p \quad (8.9)$$

which is denoted in vector-form as:

$$\begin{bmatrix} \frac{\partial J(a_N)}{\partial a_N^*(1)} \\ \frac{\partial J(a_N)}{\partial a_N^*(2)} \\ \vdots \\ \frac{\partial J(a_N)}{\partial a_N^*(p)} \end{bmatrix} = 0 \quad (8.10)$$



It is remarked that the derivative is taken in an optimum value and since the cost function is quadratic, there is only one optimum. We will denote this optimum by  $\hat{a}_N$ .

To find this optimum we recall the expression of  $J(a_N)$  in (8.5) and the error signal  $e(n; a_N)$  from (8.2), repeated here as:

$$\begin{aligned} J(a_N) &= \frac{1}{N} \sum_{n=p}^{N-1} |e(n; a_N)|^2 \\ &= \frac{1}{N} \sum_{n=p}^{N-1} |x(n) + a_N(1)x(n-1) + \cdots + a_N(p)x(n-p)|^2 \end{aligned}$$

Then the left hand side in (8.9) becomes subsequently for each value of  $k$ :

$$\begin{aligned} \frac{\partial J(a_N)}{\partial a_N^*(k)} &= \frac{1}{N} \sum_{n=p}^{N-1} \frac{\partial e(n; a_N) e^*(n; a_N)}{\partial a_N^*(k)} \\ &= \frac{1}{N} \sum_{n=p}^{N-1} e(n; a_N) \frac{\partial e^*(n; a_N)}{\partial a^*(k)} \\ &= \frac{1}{N} \sum_{n=p}^{N-1} e(n; a_N) x^*(n-k) \end{aligned}$$

Therefore (8.10) leads to:

$$\boxed{\frac{1}{N} \sum_{n=p}^{N-1} \begin{bmatrix} x^*(n-1) \\ x^*(n-2) \\ \vdots \\ x^*(n-p) \end{bmatrix} e(n; \hat{a}_N) = 0} \quad (8.11)$$

Using the matrices defined in (8.6), and recalling the definition of the *Hermitian transpose* of  $X_N$  as:

$$X_N^H = \begin{bmatrix} x^*(p-1) & x^*(p) & \cdots & x^*(N-2) \\ x^*(p-2) & x^*(p-1) & & x^*(N-3) \\ \vdots & & \ddots & \vdots \\ x^*(0) & x^*(1) & \cdots & x^*(N-p-1) \end{bmatrix}$$

we can write (8.11) explicitly as:

$$\begin{bmatrix} x^*(p-1) & x^*(p) & \cdots & x^*(N-2) \\ x^*(p-2) & x^*(p-1) & & x^*(N-3) \\ \vdots & & \ddots & \vdots \\ x^*(0) & x^*(1) & \cdots & x^*(N-p-1) \end{bmatrix} \begin{bmatrix} e(p; \hat{a}_N) \\ e(p+1; \hat{a}_N) \\ \vdots \\ e(N-1; \hat{a}_N) \end{bmatrix} = 0 \quad (8.12)$$

These equations constitute key equations in characterizing the solution of least squares problems and are indicated as the *Orthogonality Condition*. This will be analyzed and illustrated further in Section 8.3.2.

To derive an explicit form for the solution of this *orthogonality* condition, we rewrite the error signal  $e(n; a_N)$  in (8.2) as,

$$e(n; a_N) = [x(n-1) \quad x(n-2) \quad \cdots \quad x(n-p)] \begin{bmatrix} a_N(1) \\ a_N(2) \\ \vdots \\ a_N(p) \end{bmatrix} + x(n)$$

Then the "optimal" solution  $\hat{a}_N$  to AR parameter estimation problem as given by condition (8.11) can be denoted as:

$$\begin{aligned} \frac{1}{N} \left( \sum_{n=p}^{N-1} \begin{bmatrix} x^*(n-1) \\ x^*(n-2) \\ \vdots \\ x^*(n-p) \end{bmatrix} [x(n-1) \quad x(n-2) \quad \cdots \quad x(n-p)] \right) \begin{bmatrix} \hat{a}_N(1) \\ \hat{a}_N(2) \\ \vdots \\ \hat{a}_N(p) \end{bmatrix} + \\ \frac{1}{N} \sum_{n=p}^{N-1} \begin{bmatrix} x^*(n-1) \\ x^*(n-2) \\ \vdots \\ x^*(n-p) \end{bmatrix} x(n) = 0 \end{aligned} \quad (8.13)$$

These equations are the *normal equations* that characterizes the solution to the AR LLSQ problem (8.5). They can be written explicitly as:

$$\begin{bmatrix} \frac{1}{N} \sum_{n=p}^{N-1} |x(n-1)|^2 & \cdots & \frac{1}{N} \sum_{n=p}^{N-1} x(n-p)x^*(n-1) \\ \vdots & \ddots & \vdots \\ \frac{1}{N} \sum_{n=p}^{N-1} x(n-1)x^*(n-p) & \cdots & \frac{1}{N} \sum_{n=p}^{N-1} |x(n-p)|^2 \end{bmatrix} \begin{bmatrix} \hat{a}_N(1) \\ \hat{a}_N(2) \\ \vdots \\ \hat{a}_N(p) \end{bmatrix} + \begin{bmatrix} \frac{1}{N} \sum_{n=p}^{N-1} x(n)x^*(n-1) \\ \vdots \\ \frac{1}{N} \sum_{n=p}^{N-1} x(n)x^*(n-p) \end{bmatrix} = 0$$

Using again the matrices defined in (8.6), and again recalling the definition of the *Hermitian transpose* of  $X_N$ , this solution can be *compactly* denoted as:

$$\boxed{\left( \frac{1}{N} X_N^H X_N \right) \hat{a}_N + \frac{1}{N} X_N^H y_N = 0} \quad (8.14)$$

**Remark 8.1** (The Normal Equations for the Linear Regression Problem). *When abstracting the matrix  $X_n$  and the vectors  $y_N$  and  $a_N$  to other problems, such as illustrated in Example 8.3, following the same way of derivation the solution to a general linear regression problem is also characterized via the above Normal Equations (8.14).*

The Normal Equations define the estimate  $\hat{a}_N$  implicitly and an explicit definition becomes possible when the matrix  $\left( \frac{1}{N} X_N^H X_N \right)$  is invertible. In that case, the estimate  $\hat{a}_N$  is explicitly given as:

$$\hat{a}_N = - \left( \frac{1}{N} X_N^H X_N \right)^{-1} \frac{1}{N} X_N^H y_N \quad (8.15)$$

### 8.3.2 The Orthogonality Condition

The condition (8.11) or its explicit form (8.12) are called the *Orthogonality Condition* of the LLSQ problem for estimating the parameters of an AR model. Orthogonality conditions characterize the solution of many quadratic optimization problems. For that reason we are going to investigate this condition (for the case of AR parameters estimation) more closely. This will lead to a *geometric interpretation* that leads to an alternative derivation of the Normal Equations (8.14).

Let the matrix  $X_N$  defined in (8.6) be denoted as a collection of  $N - p$  dimensional vectors,

$$X_N = [x_{N,1} \quad x_{N,2} \quad \cdots \quad x_{N,p}]$$

with  $x_{N,j} \in \mathbb{C}^{N-p}$  for  $j = 1 : p$ . Then the error  $e_N$  can be written as,

$$e_N = y_N + a_N(1)x_N(1) + \cdots + a_N(p)x_N(p)$$

With these definitions the condition (8.12) can be written as,

$$x_{N,j}^H e_N = 0 \quad \text{for } j = 1 : p$$

This condition means that the vector  $e_N$  is *orthogonal* to all vectors  $x_{N,j}$  for  $j = 1 : p$ . This orthogonality leads to an alternative derivation of the least squares solution as illustrated in the following Example.

---

#### Example 8.4 (Illustration Orthogonality Condition)

Consider the AR LLSQ problem (8.5) for  $p = 2$  and  $N = 5$ , then for this case the data of the LLSQ problem given in (8.6) reads:

$$x_{5,1} = \begin{bmatrix} x(1) \\ x(2) \\ x(3) \end{bmatrix} \quad x_{5,2} = \begin{bmatrix} x(0) \\ x(1) \\ x(2) \end{bmatrix} \quad y_5 = \begin{bmatrix} x(2) \\ x(3) \\ x(4) \end{bmatrix} \quad e_5 = \begin{bmatrix} e(2; a_5) \\ e(3; a_5) \\ e(4; a_5) \end{bmatrix}$$

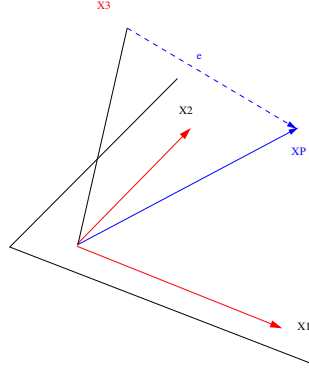
and the AR LLSQ problem can be denoted as:

$$\min_{a_5} \|y_5 - \underbrace{(-a_5(1)x_{5,1} - a_5(2)x_{5,2})}_{\text{linear combination}}\|_2^2 \quad (8.16)$$

The geometric information now is the following. The underbraced quantity represents a linear combination of the two vectors  $x_{5,1}$  and  $x_{5,2}$ . When these vectors are independent as displayed in Figure 8.3 then they span a *plane* and by varying the coefficients  $-a_5(1)$  and  $-a_5(2)$  the linear combination can be *any vector* in that 2-dim plane. The difference  $y_5 - (-a_5(1)x_{5,1} - a_5(2)x_{5,2})$  is depicted in Figure 8.3 by the dashed line for an arbitrary pair of coefficients  $-a_5(1), -a_5(2)$ . Since the LLSQ criterium  $J(a_5)$  is the length of the vector  $e_5$ , the orthogonality conditions directly follows when requiring this length to be *minimal*. That is that case when  $e_5$  is orthogonal to the plane spanned by the vectors  $x_{5,1}$  and  $x_{5,2}$ .

---

Example 8.4 illustrates that the geometrical interpretation of the minimizing of the LLSQ cost function corresponds to the minimization of the distance of a given vector  $y_N$  of dependent variables to a plane spanned by the vectors  $x_{N,j}$  for  $j = 1 : p$  (constructed from the independent variables). The condition for minimality of the distance is simply the orthogonality condition (8.11).



**Figure 8.3:** Geometric Illustration of the Orthogonality condition for  $p = 2$  and  $N = 5$  in Example 8.4.

## 8.4 The accuracy of the LLSQ estimated parameters

In this section we consider the Linear regression problem considered in Section 8.2.2. To analyse the accuracy a simple statistical scenario is considered. Let the vector of true parameter values  $a(i)$  be stored in the vector  $a$  as  $[a(1) \ a(2) \ \cdots \ a(p)]^T$ , then we consider the following specific form of the data equation (8.8):

$$y_N + X_N a = v_N \quad (8.17)$$

for the matrix  $X_N$  deterministic and  $v_N$  a random vector with zero-mean entries and covariance matrix:

$$E[v_N v_N^H] = \sigma_v^2 I_N \quad (8.18)$$

where  $I_N$  is the  $N \times N$  identity matrix. The vector of true parameters  $a$  is also deterministic.

This statistical scenario makes the vector  $y_N$  a vector of random variables and as a consequence, the LLSQ estimate  $\hat{a}_N$  as given by (8.15) also is a random vector. This is why on page 151 we referred to it as an estimate.

The quality of an estimate can be characterized by its probability density function. However as an alternative, in Chapter 5, we could restrict the characterization of a random variable to its Ensemble averages. Two classically used ensemble averages are the mean and the covariance matrix of a random vector. These are denoted as follows:

1. *Mean:*  $E[\hat{a}_N]$
2. *Covariance Matrix:*  $E[(\hat{a}_N - E[\hat{a}_N])(\hat{a}_N - E[\hat{a}_N])^H]$

For the simple statistical scenario we first evaluate the mean  $E[\hat{a}_N]$ . For that purpose we substitute the expression for  $y_N = -X_N a + v_N$  into the LLSQ solution (8.15) to obtain:

$$\hat{a}_N = \left( \frac{1}{N} X_N^H X_N \right)^{-1} \frac{1}{N} X_N^H X_N a - \left( \frac{1}{N} X_N^H X_N \right)^{-1} \frac{1}{N} X_N^H v_N \quad (8.19)$$

Since  $X_N$  was assumed to be deterministic and  $E[v_N] = 0$ , the mean  $E[\hat{a}_N]$  is:

$$E[\hat{a}_N] = a - \left( \frac{1}{N} X_N^H X_N \right)^{-1} \frac{1}{N} X_N^H E[v_N] = a \quad (8.20)$$

Therefore the mean  $E[\hat{a}_N]$  equals  $a$ , the true parameter vector. Such an estimate is called *unbiased*.

Second we evaluate under the same conditions the covariance matrix  $E[(\hat{a}_N - E[\hat{a}_N])(\hat{a}_N - E[\hat{a}_N])^H]$ . Again using the expression for  $y_N = -X_N a + v_N$ , the expression for  $\hat{a}_N$  in (8.19) and the unbiased property of the estimate  $E[\hat{a}_N]$  the covariance matrix becomes:

$$\begin{aligned} E[(\hat{a}_N - E[\hat{a}_N])(\hat{a}_N - E[\hat{a}_N])^H] &= E[(\hat{a}_N - a)(\hat{a}_N - a)^H] \\ &= E\left[\left(\frac{1}{N} X_N^H X_N\right)^{-1} \frac{1}{N} X_N^H v_N v_N^H \frac{1}{N} X_N \left(\frac{1}{N} X_N^H X_N\right)^{-1}\right] \\ &= \left(\frac{1}{N} X_N^H X_N\right)^{-1} \frac{1}{N} X_N^H E[v_N v_N^H] \frac{1}{N} X_N \left(\frac{1}{N} X_N^H X_N\right)^{-1} \end{aligned}$$

For the assumed covariance matrix of the random vector  $v_N$  given in (8.18), the covariance matrix of the estimated parameters  $\hat{a}_N$  becomes:

$$E[(\hat{a}_N - E[\hat{a}_N])(\hat{a}_N - E[\hat{a}_N])^H] = \sigma_v^2 \left( X_N^H X_N \right)^{-1} \quad (8.21)$$

**Remark 8.2** (Application to AR parameter estimation problem). *For the AR model (7.11) the data equation in (8.17) with the definition of  $y_N$  and  $X_N$  as in (8.6) can still be used if we define the vector  $v_N$  as  $[v(p) \ v(p+1) \ \cdots \ v(N-1)]$ . Indeed in this case the vector  $v_N$  satisfies the condition (8.18). The matrix  $X_N$  is also assumed to be deterministic since it is constructed from a single given realization of finite samples of the stochastic process as denoted in (8.1). The observation  $y_N$  in the data equation becomes random due to the presence of the (unknown) vector  $v_N$ . Based on these observations, we can apply that simple statistical framework to derive the mean and the covariance matrix of the LLSQ AR estimated parameters as given resp. in (8.20) and (8.21).*

---

## References

- [1] K. Hinnen, *Data-Driven Optimal Control for Adaptive Optics*, PhD thesis. Delft University of Technology, 2007.
- [2] M.A. van Dam, D. Le Mignant, B.A. Macintosh, *Performance of Keck observatory adaptive optics system*. Applied Optics 43 (29), 5458–5467, 2004.
- [3] M. Verhaegen, G. Vdovin and O. Soloviev, *Control for High Resolution Imaging*. Lecture Notes, Delft University of Technology (Course sc4045), 2015.
- [4] M. Verhaegen and V. Verdult, *Filtering and System Identification: A Least Squares Approach*. Cambridge: Cambridge University Press, 2007. (Used in Course sc4040).

- [5] D. Clay, *Linear Algebra and its Applications*, 4th-Ed., Addison-Wesley, 2012.
- [6] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.

## Exercises

**Exercise 8.1** Consider the following AR model for  $p = 1$  and  $v(n)$  ZMWN with variance equal to 1:

$$x(n) - ax(n-1) = v(n) \quad |a| < 1$$

then show that

$$E[x(n-1)v^*(n)] = 0$$

[In the same way you can show:

$$E[x(n-\alpha)v^*(n)] = 0 \quad \alpha \geq 1 \quad ]$$

**Exercise 8.2** Consider the LLSQ problem outlined in Example 8.3, repeated explicitly here:

$$\min_{R_N, u_{0,N}} \left\| \begin{bmatrix} i(0) & 1 \\ i(1) & 1 \\ \vdots & \vdots \\ i(N-1) & 1 \end{bmatrix} \begin{bmatrix} R_N \\ u_{0,N} \end{bmatrix} - \begin{bmatrix} u_m(0) \\ u_m(1) \\ \vdots \\ u_m(N-1) \end{bmatrix} \right\|_2^2$$

Then show that the optimal solution to this LLSQ problem is determined by the following Normal Equations:

$$\begin{bmatrix} \sum_{n=0}^{N-1} i(n)^2 & \sum_{n=0}^{N-1} i(n) \\ \sum_{n=0}^{N-1} i(n) & N \end{bmatrix} \begin{bmatrix} \hat{R}_N \\ \hat{u}_{0,N} \end{bmatrix} = \begin{bmatrix} \sum_{n=0}^{N-1} i(n)u_m(n) \\ \sum_{n=0}^{N-1} u_m(n) \end{bmatrix}$$

**Exercise 8.3** Consider the following least squares problem:

$$\min_a \|Xa - y\|_2^2$$

with  $X \in \mathbb{R}^{N \times m}$  and  $y \in \mathbb{R}^N$ ,  $a \in \mathbb{R}^m$  given resp. as:

$$X = \begin{bmatrix} x_1(1) & x_2(1) & \cdots & x_m(1) \\ x_1(2) & x_2(2) & \cdots & x_m(2) \\ \vdots & & \ddots & \vdots \\ x_1(N) & x_2(N) & \cdots & x_m(N) \end{bmatrix} \quad y = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{bmatrix} \quad a = \begin{bmatrix} a(1) \\ a(2) \\ \vdots \\ a(m) \end{bmatrix}$$

With these matrices we consider the least squares problem:

$$\min_a \|Xa - y\|_2^2 = \min_a J(a)$$

- (a) Then using the condition that characterizes the optimum to the considered least squares problem:

$$\begin{bmatrix} \frac{\partial J(a)}{\partial a(1)} \\ \frac{\partial J(a)}{\partial a(2)} \\ \vdots \\ \frac{\partial J(a)}{\partial a(p)} \end{bmatrix} = 0 \quad (8.22)$$

derive the matrix expression of the orthogonality condition (8.12).

- (b) If we denote the matrix  $X$  in terms of its column vectors as:

$$X = [x_1 \quad x_2 \quad \cdots \quad x_m]$$

with each  $x_i$  (for  $i = 1 : m$ ) a column vector in  $\mathbb{R}^N$ , and if we denote the least squares estimate by  $\hat{a}$  and the error vector as  $e = X\hat{a} - y$ , then show that the *Orthogonality condition* (8.12) can explicitly be written as:

$$X^T e = 0$$

**Exercise 8.4** When the relationship between two signals  $x(n)$  and  $y(n)$  is given by the following law:

$$y(n) = \alpha e^{-\beta x(n)}$$

then formulate a Linear Regression problem to estimate the parameters  $\alpha$  and  $\beta$  assuming that the signals  $x(n)$  and  $y(n)$  are exactly known for  $n = 0 : N - 1$ .

**Exercise 8.5** When the stochastic process  $x(n)$  in the estimation of the AR parameters is Auto-correlation ergodic, such that the following limits hold:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=p}^{N-1} x(n - \alpha) x^*(n - \beta) = r_x(\beta - \alpha)$$

for  $r_x(\beta - \alpha) = E[x(n - \alpha) x^*(n - \beta)]$ , then show

- (a) that the normal equations that characterizes the solution to the AR LLSQ problem (8.5) taking the limit for  $N \rightarrow \infty$  becomes equal to,

$$\begin{bmatrix} r_x(0) & r_x^*(1) & \cdots & r_x^*(p-1) \\ r_x(1) & r_x(0) & \cdots & r_x^*(p-2) \\ \vdots & & \ddots & \\ r_x(p-1) & r_x(p-2) & \cdots & r_x(0) \end{bmatrix} \begin{bmatrix} \hat{a}(1) \\ \hat{a}(2) \\ \vdots \\ \hat{a}(p) \end{bmatrix} + \begin{bmatrix} r_x(1) \\ r_x(2) \\ \vdots \\ r_x(p) \end{bmatrix} = 0$$

Where the estimated parameters  $\hat{a}_N(i)$  in the limit are denoted by  $\hat{a}$ .

- (b) Show that the solution  $\hat{a}(i)$  in the limit  $N \rightarrow \infty$  corresponds to the set of equations (7.16) to determine the AR parameters in the inverse AR problem.

**Exercise 8.6** Consider the estimation of the value of Resistor as in Example 8.3 with measurement model again noted as:

$$u_m(n) = Ri(n) + u_0 + e(n)$$

Assume that  $e(n)$  is ZMWN with variance  $\sigma_e^2$ . Consider the unknown quantities  $R, u_0$  and the current sequence  $i(n)$  to be deterministic. We have  $N$  observations for  $n = 0 : N - 1$ , then

- (a) determine the LLSQ estimate of the following LLSQ problem:

$$\min_{R_N} \sum_{n=0}^{N-1} \left| u_m(n) - R_N i(n) \right|^2$$

Let this estimate be denoted as  $\hat{R}_N$ .

- (b) Determine the mean and covariance matrix of the estimate  $\hat{R}_N$ .  
(c) Is the estimate  $\hat{R}_N$  unbiased?

**Exercise 8.7** Consider the current  $I$  in a resistor, where the resistance is a random variable  $R$ , uniformly distributed on  $[9, 11]$ . The voltage  $V$  over the resistor is a Gaussian distributed random variable with expected value  $\mu_V = 9$  and standard deviation  $\sigma_V = 2$ . Assume that the random variables  $R$  and  $V$  are statistically independent. Ohm's law tells us that:

$$I = \frac{V}{R}.$$

- (a) Compute the expected value  $\mu_I$  and the variance  $\sigma_I^2$  of the current  $I$ .

Hint:  $\int \frac{1}{x} dx = \ln |x|$ .

- (b) Compute the correlation  $r_{IV}$  and the correlation coefficient  $\rho_{IV}$  between the current  $I$  and the voltage  $V$ .

Hint:  $r_{IV} = E[IV]$  and  $\rho_{IV} = \frac{\sigma_{IV}}{\sigma_I \sigma_V}$ , with  $\sigma_{IV} = E[(I - \mu_I)(V - \mu_V)]$  the covariance between  $I$  and  $V$ .

- (c) Suppose we have measured the voltage  $V$  and want to estimate the current  $I$  based on this measurement. The linear mean-squared error estimator of  $I$  based on  $V$  can then be written as:

$$\hat{I} = a + bV.$$

Determine the coefficients  $a$  and  $b$  (in terms of  $\mu_V, \mu_I, \rho_{IV}, \sigma_I, \sigma_V$ ).

**Exercise 8.8** Suppose we have  $n$  observations  $y_1, \dots, y_n$ , described by:

$$y_i = \alpha x_i + \beta + e_i, \quad i = 1, \dots, n.$$

where the values  $x_i$  are known exactly. The random disturbances  $e_i$  are independent and have an expected value of zero and a variance of  $\sigma^2$ . Show that the least squares estimators of the parameters  $\alpha$  and  $\beta$  are uncorrelated if the mean of all  $x_i$  equals zero, i.e.  $\frac{1}{n} \sum_{i=1}^n x_i = 0$ .

What is the variance of the estimators under this condition?



**Exercise 8.9** A satellite moves along a straight line for a certain time and with a constant velocity  $\beta$ . The initial position of the satellite is  $\alpha$ . The observed position of the satellite is given by:

$$y_t = \alpha + \beta t + e_t, \quad t = 1, 2, \dots, n$$

where the noise  $e_t$  is a random variable with expected value 0 and variance 1.

Suppose the observations are given by

| t  | $y_t$ |
|----|-------|
| 1  | 2     |
| 2  | 2     |
| 3  | 3     |
| 4  | 4     |
| 5  | 4     |
| 6  | 8     |
| 7  | 9     |
| 8  | 10    |
| 9  | 11    |
| 10 | 12    |

- Compute the (unbiased) least squares estimators  $\hat{\alpha}$  and  $\hat{\beta}$  of  $\alpha$  and  $\beta$ .
- Compute the variance of the least squares estimators  $\hat{\alpha}$  and  $\hat{\beta}$  of the parameters  $\alpha$  and  $\beta$ .
- Are the estimators  $\hat{\alpha}$  and  $\hat{\beta}$  uncorrelated? Please motivate your answer.

Hint: if necessary, use MATLAB for solving this exercise.

**Exercise 8.10** The (theoretical) relationship between the light intensity in an optical fiber and the distance to the light source is given by:

$$I(x) = I_0 e^{-\alpha x}, \quad (8.23)$$

with  $I(x)$  the light intensity at a distance  $x$  of the light source,  $\alpha$  the absorption coefficient and  $I_0$  the light intensity of the source. Suppose we measure the light intensity on  $n$  different points  $\{x_i\}_{i=1,\dots,n}$  in the optical fiber, resulting in the measurements  $\{I(x_i)\}_{i=1,\dots,n}$ . The positions  $\{x_i\}_{i=1,\dots,n}$  are assumed to be exactly known (measured noise-free), while the measurements  $\{I(x_i)\}_{i=1,\dots,n}$  are noise-disturbed and, therefore, are to be considered as realizations of the random variables  $\{I(x_i)\}_{i=1,\dots,n}$ .

- We first consider the situation that the parameters  $I_0$  and  $\alpha$  are both unknown and must be estimated from the measurements. Based on the theoretical relationship (8.23), it follows:

$$\log I(x) = \log I_0 - \alpha x,$$

The following least-squares estimator of  $\theta = \begin{pmatrix} \log I_0 \\ \alpha \end{pmatrix}$  can be designed:

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^n (\log I(x_i) + \alpha x_i - \log I_0)^2$$

What are the requirements on the noise disturbances on the measurements, for this estimator to be unbiased? Motivate your answer.

- (b) Now, we assume that the parameter  $\alpha$  is known, and that we want to estimate  $I_0$  based on the measurements  $\{x_i, I(x_i)\}_{i=1, \dots, n}$ . Furthermore, the measurements  $\{I(x_i)\}_{i=1, \dots, n}$  can be modeled as independent, exponentially distributed random variables with expected values  $\mathcal{E}[I(x_i)] = I_0 \exp(-\alpha x_i)$ .

Give a (closed form) expression for the (unweighted) least-squares estimator of  $I_0$  and determine the expected value and variance of this estimator.

# Chapter 9

## Optimal Filtering

**After studying this chapter you can:**

- formulate the applications of optimal filtering, i.e. denoising a noisy signal, deconvolution, prediction and active noise cancellation that were presented in Section 1.5 as an optimal filter design problem
- derive the Wiener-Hopf equations for determining the optimal parameters of a filter with a so-called FIR transfer function  $W(z)$  by minimizing the variance of the error signal between the signal of interest and the output of  $W(z)$  and in addition determine the variance of this error signal for the optimal filter parameters
- derive the Wiener-Hopf equations for determining the parameters of a filter with a so-called mixed causal, anti-causal IIR transfer function  $W(z)$  by minimizing the variance of the error signal between the signal of interest and the output of  $W(z)$  and in addition determine the variance of this error signal for the optimal filter parameters
- derive the Wiener-Hopf equations for determining the parameters of a filter with a causal IIR transfer function  $W(z)$  by minimizing the variance of the error signal between the signal of interest and the output of  $W(z)$
- in all the above optimization problems know what information on the data is needed.

---

## 9.1 Introduction

Looking back on what has been learned so far in the previous Chapters, it may be concluded that the reader is introduced in two parts of the field of stochastic processes.

In the first part (Chapters 5 and 6) knowledge has been provided on how to describe mainly WSS stochastic processes via their Auto-Cross correlation functions or Power- Cross Spectra. Subsequently we learned to update these statistical descriptions when filtering the given stochastic process with an LTI filter.

In the second part (Chapters 7 and 8) we looked at important inverse problems to retrieve models when only statistical information is available about the stochastic process or only a finite number of samples of a single realization of that stochastic process.

In this Chapter we explore a third and final part. More precisely we use the insights built up in the previous two parts to ‘polish’ or ‘shape’ stochastic processes. Such polishing might mean to reduce the effect of ‘unwanted’ noise that has corrupted the measurement of a realization of the stochastic process. In this introductory context we formulate this ‘polishing problem’ as a filter design problem, with the filter filtering the given signal such that the output approximates the desired signal. This filter design problem is a so-called *synthesis problem*. Such problems are contrary to the analysis problems of Chapter 6 where the filter was given.

As highlighted in Chapter 1 four different applications of designing a “shaping filter” will be considered:

1. Denoising
2. Prediction
3. Deconvolution
4. Active Noise Cancellation

Though these problems are different in nature, as in [1] we will first address a generic filter design problem and then specialize to a number of these applications. This will be done as in [1] first for the filter to be designed having a Finite Impulse Response (FIR) structure and then we will consider the required filter to have an Infinite Impulse Response (IIR) structure. The FIR case is treated in Section 9.3 and the IIR case in Section 9.4. When formulating the design in a minimum variance framework the resulting filters are so-called *Wiener filters*.

---

## 9.2 Generic Filtering Problem

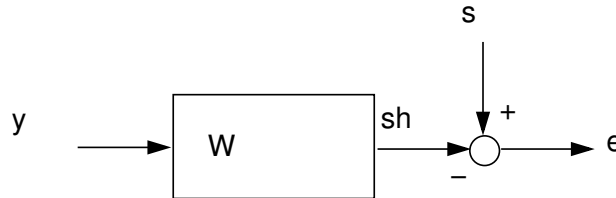
Explained in generic terms the goal we are interested in this chapter is to estimate a stochastic process  $d(n)$  that is assumed not to be measurable, from another stochastic process  $x(n)$  that is measurable. The estimation will be formulated as

the problem of designing an LTI filter with transfer function  $W(z)$ , that filters the signal  $x(n)$  such that its output is an estimate of  $d(n)$ , denoted as  $\hat{d}(n)$ , in some optimal way. The optimality here is expressed by a Minimum Variance Cost Function. As we may understand intuitively from this problem formulation, the degree of success in retrieving one signal from another will of course depend on the relation between both signals. In the context of these course notes, this relationship is expressed in terms of the (Auto- and Cross Correction functions or Spectra) between the signals of interest. The standing assumption in this Chapter will be that all signals are considered to be WSS.

To define this generic problem, consider Figure 9.1. Here the notation  $W(z; w(j))$  is used to explicitly indicate that the filter depends on (tunable) parameters  $w(j)$ , where the index  $j \in \mathcal{S} \subset \mathbb{Z}$  as specified by the specific filter at hand. A consequence is that the resulting filtered signals  $\hat{d}(n; w(j))$  and  $e(n; w(j))$  also depend on the filter parameters. With this notation the goal is then to design (tune) these filter parameters  $w(j)$  in  $W(z; w(j))$  in order to minimize the variance of the error signal  $e(n; w(j))$ :

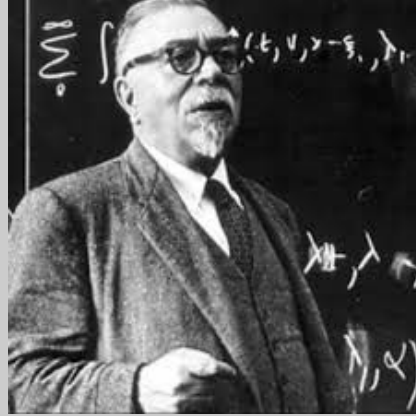
$$E[|e(n)|^2] = E[|d(n) - \hat{d}(n)|^2] \quad (9.1)$$

This is referred to as *Minimum Variance Filter Design*. In order to make the Filter



**Figure 9.1:** Block-schematic representation of the generic filter design problem to estimate signal  $\hat{d}(n; w(j))$  from the signal  $x(n)$  by designing the filter with transfer function  $W(z; w(j))$  depending on the filter parameters  $w(j)$  for  $j \in \mathcal{S} \subset \mathbb{Z}$ .

Design problem tractable, the parametrization of the filter  $W(z; w(j))$  is important. We will consider two parametrizations in the next two sections. The first in Section 9.3, where  $W(z)$  is parametrized as an FIR transfer function and the second in Section 9.4, where  $W(z)$  is parametrized as an IIR transfer function. The optimal resulting filters are called *Wiener filters* referring to the pioneer in this field Norbert Wiener. He was an American mathematician who coined the term of science of *cybernetics*. This is the field of science that is dealing with studying and optimizing the control, operation and communication between living organisms, autonomous or automatized machines and their infrastructure.



(Adapted from [https://en.wikipedia.org/wiki/Norbert\\_Wiener](https://en.wikipedia.org/wiki/Norbert_Wiener))

Norbert Wiener was borne in 1894 as the first child of Jewish parents of Polish and German origin. Norbert was educated by his father at home until 1903. Norbert described his father as calm and patient, unless he (Norbert) failed to give a correct answer, at which his father would lose his temper.

After graduating from High School in 1906 at the age of 11, Wiener entered Tufts College. He was awarded a BA in mathematics in 1909 at the age of 14, whereupon he received a Ph.D. in 1912, when he was merely 17 years old. His supervisor was the mathematician Karl Schmidt.

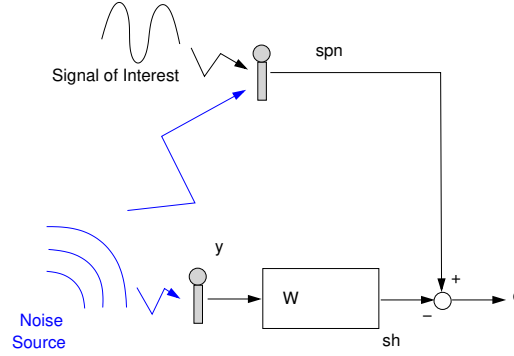
With the generic Block-scheme in Figure 9.1 three applications of interest can be defined.

**Denoising:** In that case we observe a noisy variant of a desired signal  $d(n)$  and we assume that this signal is disturbed by additive noise  $v(n)$ . The noisy signal available for processing is the signal  $x(n)$ :

$$x(n) = d(n) + v(n)$$

**Prediction:** Here we would like to predict the observation  $d(n)$  into the future. For the  $\alpha$ -step ahead prediction case with  $\alpha \in \mathbb{N}$  the signal  $x(n) = d(n)$  and the signal of interest is  $d(n + \alpha)$ .

**Active Noise Cancellation:** The problem was mentioned in Section 1.5.4 and the block-scheme is repeated in Figure 9.2. The goal is to estimate the signal  $v_1(n)$  from the signal  $v_2(n)$ . This estimate is then used to ‘clean’ the corrupted signal  $d(n) + v_1(n)$ . To bring this in the generic filtering schematic of Figure 9.1, there is an overlap in the use of the symbol  $d(n)$ . However it should be clear from the context what is meant. For active noise cancellation in the context of Figure 9.1 we then would take  $d(n) = v_1(n)$  and  $x(n) = v_2(n)$ .



**Figure 9.2:** The signal  $d(n) + v_1(n)$  is the noise corrupted signal. The filter  $W(z)$  filters the signal  $v_2(n)$  to produce the estimate  $\hat{v}_1(n)$ .

Apart from the mathematical formulation and solution of an optimal filter design (or synthesis) problem, it is also important to analyse what the practical requirements are to calculate the solution.

## 9.3 Minimum Variance FIR Wiener Filter

### 9.3.1 The Generic Problem Formulation

For the FIR Wiener filter design problem the filter  $W(z; w(j))$  in the generic Block-scheme in Figure 9.1 has the following form:

$$W(z; w(j)) = \sum_{j=0}^{m-1} w(j)z^{-j} \quad (9.2)$$

This results in the following difference equation relating the input  $x(n)$  to the output  $\hat{d}(n; w(j))$ :

$$\hat{d}(n; w(j)) = w(0)x(n) + w(1)x(n-1) + \dots + w(m-1)x(n-m+1) = \sum_{k=0}^{m-1} w(k)x(n-k) \quad (9.3)$$

If we compare the summation in this equation with that in the general convolution expression given in (2.11), we see that the summation in (9.3) is finite. Therefore the impulse response is finite, explaining the name FIR model.

For this model the *Minimum Variance FIR Wiener Problem* can now be stated as:

$$\min_{w(i)} J(w(j)) \quad \text{for } J(w(j)) = E[|e(n; w(j))|^2] \quad (9.4)$$

The integer  $m$  is called the *order* of the FIR model.

### 9.3.2 Solution to the Generic Problem

The solution to the Generic FIR Wiener Problem is summarized in the following Theorem.

**Theorem 9.1** (FIR Wiener Filter). *Let the conditions as stipulated in the generic problem formulation given in Section 9.3.1 hold, let the stochastic processes  $x(n)$  and  $d(n)$  be jointly zero mean and WSS, let the Auto-correlation function of  $d(n)$  be denoted as  $r_d(k)$ , and let the following vector quantities (in bold) be defined:*

$$\mathbf{x}^*(\mathbf{n}) = \begin{bmatrix} x^*(n) \\ \vdots \\ x^*(n-m+1) \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w(0) \\ \vdots \\ w(m-1) \end{bmatrix} \quad (9.5)$$

and the following covariance matrix and cross-correlation vector:

$$\mathbf{R}_x = E[\mathbf{x}^*(\mathbf{n})\mathbf{x}(\mathbf{n})^T] > 0 \quad \mathbf{r}_{dx} = E[d(n)\mathbf{x}^*(\mathbf{n})] \quad (9.6)$$

then the solution to the Minimum Variance FIR Wiener Problem (9.4) is given in terms of the optimal filter coefficient vector  $\hat{\mathbf{w}}$  and value of the criterion at this optimum as:

$$\mathbf{R}_x \hat{\mathbf{w}} = \mathbf{r}_{dx} \quad (\text{The Wiener-Hopf equations}) \quad (9.7)$$

$$J(\hat{\mathbf{w}}) = r_d(0) - \mathbf{r}_{dx}^H \mathbf{R}_x^{-1} \mathbf{r}_{dx} \quad (9.8)$$

with the notation  $(.)^H$  used to indicate the Hermitian transpose.

*Proof.* As for the AR parameter estimation problem treated in Section 8.2.1 the determination of the optimal filter coefficient vector  $\hat{\mathbf{w}}$  that minimizes the variance of the error signal  $E[|e(n; \hat{\mathbf{w}})|^2]$  may rely on the derivation of the *orthogonality condition* for the problem at hand.

To derive this condition first we express the criterion as:

$$\begin{aligned} J(\mathbf{w}) &= E[|e(n; \mathbf{w})|^2] \\ &= E[|d(n) - \hat{d}(n; \mathbf{w})|^2] \\ &= E\left[\left|d(n) - \sum_{\ell=0}^{m-1} w(\ell)x(n-\ell)\right|^2\right] \end{aligned} \quad (9.9)$$

and subsequently consider the following derivative (now given in vector form):

$$\begin{bmatrix} \frac{\partial J(\mathbf{w})}{\partial w^*(0)} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w^*(m-1)} \end{bmatrix}$$

Using the expression for  $J(\mathbf{w})$  given in equation (9.9) and the fact that the derivative is taken with respect to (w.r.t.) the parameters  $w(j)$ , this derivative becomes (as explained for the two parameter case in Section 2.4):

$$\begin{bmatrix} E\left[\frac{\partial e(n; \mathbf{w})e^*(n; \mathbf{w})}{\partial w^*(0)}\right] \\ \vdots \\ E\left[\frac{\partial e(n; \mathbf{w})e^*(n; \mathbf{w})}{\partial w^*(m-1)}\right] \end{bmatrix} = \begin{bmatrix} E\left[e(n; \mathbf{w})\frac{\partial e^*(n; \mathbf{w})}{\partial w^*(0)}\right] \\ \vdots \\ E\left[e(n; \mathbf{w})\frac{\partial e^*(n; \mathbf{w})}{\partial w^*(m-1)}\right] \end{bmatrix}$$



Now consider the expression for the error signal as given in the last equation in (9.9), the necessary and sufficient conditions for optimality of the quadratic criterion  $J(\mathbf{w})$  leads to the *orthogonality condition* for the Minimum Variance FIR Wiener Problem:

$$E \left[ e(n; \hat{\mathbf{w}}) \begin{bmatrix} x^*(n) \\ \vdots \\ x^*(n-m+1) \end{bmatrix} \right] = 0 \quad (9.10)$$

Writing the error signal  $e(n; \hat{\mathbf{w}})$  as,

$$e(n; \hat{\mathbf{w}}) = d(n) - [x(n) \ \cdots \ x(n-m+1)] \hat{\mathbf{w}} \quad (9.11)$$

Then involving the linearity of the expectation operator  $E[\cdot]$ , we can write the orthogonality conditions (9.10) as:

$$E \left[ \begin{bmatrix} x^*(n) \\ \vdots \\ x^*(n-m+1) \end{bmatrix} d(n) \right] - E \left[ \begin{bmatrix} x^*(n) \\ \vdots \\ x^*(n-m+1) \end{bmatrix} [x(n) \ \cdots \ x(n-m+1)] \right] \hat{\mathbf{w}} = 0$$

With the definition of the covariance matrix  $\mathbf{R}_x$  and cross-correlation vector in (9.6) this is the *Wiener-Hopf equation*. Since the matrix  $\mathbf{R}_x$  is positive definite it is invertible and the solution to the Wiener-Hopf equation is given as:

$$\hat{\mathbf{w}} = \mathbf{R}_x^{-1} \mathbf{r}_{dx} \quad (9.12)$$

To find the value of the criterion  $J(\mathbf{w})$  at the optimal filter coefficients  $\mathbf{w} = \hat{\mathbf{w}}$  we write down the following series of expressions:

$$\begin{aligned} J(\hat{\mathbf{w}}) &= E[e(n; \hat{\mathbf{w}}) e^*(n; \hat{\mathbf{w}})] \\ &= E[e(n; \hat{\mathbf{w}}) d^*(n)] - E[e(n; \hat{\mathbf{w}}) [x^*(n) \ \cdots \ x^*(n-m+1)]] \begin{bmatrix} \hat{w}^*(0) \\ \vdots \\ \hat{w}^*(m-1) \end{bmatrix} \\ &= E \left[ \left( d(n) - [x(n) \ \cdots \ x(n-m+1)] \begin{bmatrix} \hat{w}(0) \\ \vdots \\ \hat{w}(m-1) \end{bmatrix} \right) d^*(n) \right] \\ &= E[d(n) d^*(n)] - E \left[ d^*(n) [x(n) \ \cdots \ x(n-m+1)] \right] \begin{bmatrix} \hat{w}(0) \\ \vdots \\ \hat{w}(m-1) \end{bmatrix} \\ &= r_d(0) - \mathbf{r}_{dx}^H \mathbf{R}_x^{-1} \mathbf{r}_{dx} \end{aligned}$$

The first equation of this series is simply the definition of  $J(\hat{\mathbf{w}})$ , the second follows from the expression of the error signal  $e^*(n; \hat{\mathbf{w}})$  that follows from its definition in (9.11). The third equation follows from the orthogonality condition (9.10) and again using the expression of the error signal  $e(n; \hat{\mathbf{w}})$  and the fact that the optimal filter coefficients are deterministic (constant). The final expression follows from the definition of the optimal solution in (9.12), the term  $r_d(0)$  and the covariance matrix and cross correlation vector defined in (9.6)  $\square$

To conclude this section, we investigate what information is needed about the stochastic process  $x(n)$  and  $d(n)$  to calculate the optimal Wiener filter via the Wiener Hopf equation. In order to see this let the Auto-correlation function of  $x(n)$  be denoted as  $r_x(k) = E[x(n)x^*(n-k)]$  and the cross-correlation function between  $d(n)$  and  $x(n)$  be denoted as  $r_{dx}(k) = E[d(n)x^*(n-k)]$ , then the Wiener-Hopf equation is denoted in terms of these functions as:

$$\begin{bmatrix} r_x(0) & r_x^*(1) & \cdots & r_x^*(m-1) \\ r_x(1) & r_x(0) & & r_x^*(m-2) \\ \vdots & & \ddots & \vdots \\ r_x(m-1) & \cdots & & r_x(0) \end{bmatrix} \begin{bmatrix} \hat{w}(0) \\ \hat{w}(1) \\ \vdots \\ \hat{w}(m-1) \end{bmatrix} = \begin{bmatrix} r_{dx}(0) \\ r_{dx}(1) \\ \vdots \\ r_{dx}(m-1) \end{bmatrix} \quad (9.13)$$

From this explicit expression of the Wiener-Hopf equations it is clear that the following samples of the Auto- and Cross correlation functions are necessary:

$$r_x(0), r_x(1), \dots, r_x(m-1) \quad r_{dx}(0), r_{dx}(1), \dots, r_{dx}(m-1)$$

This means that not only information is necessary about the measured signal **but also** information on the (cross-) correlation between the measured signal  $x(n)$  and the signal  $d(n)$  to be estimated. Whether that is a realistic assumption is studied in more detail in the following three subsections.

### 9.3.3 Application to the Denoising Problem

In the denoising problem we consider an additive perturbation or noise  $v(n)$  to the signal  $d(n)$  of interest. The additive noise  $v(n)$  is assumed also to be zero-mean WSS and satisfies,

$$E[d(n)v^*(n-k)] = 0 \quad \forall k \quad (9.14)$$

Then we specialize the signals  $x(n)$  and  $d(n)$  in the generic case pictured in Figure 9.1 as related as follows:

$$x(n) = d(n) + v(n) \quad (9.15)$$

For these specific assumptions the data necessary in the Wiener-Hopf equations:

$$r_x(k) = E[x(n)x^*(n-k)] \quad r_{dx}(k) = E[d(n)x^*(n-k)]$$

is studied in more detail.

Using the relation between  $x(n)$  and  $d(n)$  as given in (9.15) and making use of the assumption on  $v(n)$  in (9.14), we obtain resp.

$$\begin{aligned} r_x(k) &= E[(d(n) + v(n))(d^*(n-k) + v^*(n-k))] \\ &= r_d(k) + r_v(k) \quad \text{and} \end{aligned} \quad (9.16)$$

$$\begin{aligned} r_{dx}(k) &= E[d(n)(d^*(n-k) + v^*(n-k))] \\ &= r_d(k) \end{aligned} \quad (9.17)$$

Therefore to compute both the optimal Wiener filter for denoising as well as the value of the criterion at this optimum (both given in Theorem 9.1), we need for an FIR filter  $W(z)$  of order  $m$  the Auto-correlation functions of  $r_d(k)$  and  $r_v(k)$  for lags  $0 : m - 1$ .

However when only the measured signal  $x(n)$  is available, and when this stochastic process is Auto-correlation Ergodic we only have an estimate of  $r_x(k)$  that was denoted by the estimate  $\hat{r}_x(k, N)$  in (5.10) when having  $N$  samples of a realization of the stochastic process  $x(n)$  available. In order to make an estimate of the Autocorrelation function  $r_d(k)$  for lags  $k = 0 : m - 1$  we need to have information on the Autocorrelation function  $r_v(k)$  of the noise. When the noise  $v(n)$  is assumed *in addition* to be ZMWN with variance  $\sigma_v^2$  then only an estimate of that variance is necessary. This may e.g. be obtained from the instrument used to measure the signal  $x(n)$ .

### 9.3.4 Application to the Prediction Problem

In the multi-step prediction problem we aim to estimate the value of a signal of interest  $\alpha \geq 1$  steps ahead of time. Therefore when  $n$  denotes the current time instant and when the signal of interest is  $d(n)$  and we define the vector  $\mathbf{x}(\mathbf{n})$  now as:

$$\mathbf{x}(\mathbf{n}) = \begin{bmatrix} d(n) \\ d(n-1) \\ \vdots \\ d(n-m+1) \end{bmatrix}$$

we seek to minimize the following variance:

$$J(\mathbf{w}) = E[\|\mathbf{w}^H \mathbf{x}(\mathbf{n}) - d(n+\alpha)\|^2] \quad (9.18)$$

Following the proof of Theorem 9.1 the Wiener-Hopf equation that determines the optimal filter coefficients  $\hat{\mathbf{w}}$  is given as:

$$\begin{bmatrix} r_d(0) & r_d^*(1) & \cdots & r_d^*(m-1) \\ r_d(1) & r_d(0) & & r_d^*(m-2) \\ \vdots & & \ddots & \\ r_d(m-1) & r_d(m-2) & \cdots & r_d(0) \end{bmatrix} \begin{bmatrix} \hat{w}(0) \\ \hat{w}(1) \\ \vdots \\ \hat{w}(m-1) \end{bmatrix} = \begin{bmatrix} r_d(\alpha) \\ r_d(\alpha+1) \\ \vdots \\ r_d(\alpha+m-1) \end{bmatrix}$$

Therefore to predict a signal  $d(n)$   $\alpha = 1$  steps ahead in time the following values of its Auto-correlation function are necessary:

$$r_d(0), r_d(1), \cdots, r_d(m-1), r_d(m)$$

For  $\alpha = 1$  the estimate is indicated as the *one-step ahead prediction* and is commonly used in the design of optimal controllers.

We could apply the Wiener-Hopf theory to a more general prediction problem where the signal of interest can only be measured up to an additive perturbation. Exercise 9.1 calls for an analysis of this generalization.

### 9.3.5 Application to the Active Noise Cancellation Problem

The Active Noise Cancellation problem is schematically depicted in Figure 9.2. The problem is to estimate  $v_1(n)$  from the signal  $v_2(n)$ . This is the key in attenuating the additive noise on the desired signal  $d(n)$ , since the estimate  $\hat{v}_1(n)$  can then be subtracted from the measured signal  $y(n) = d(n) + v_1(n)$ . In order to address this problem, the following assumptions are made:

1. All stochastic processes are zero-mean WSS.
2. The stochastic processes  $d(n)$  and  $v_2(n)$  are uncorrelated, such that,

$$E[d(n)v_2^*(n-k)] = 0 \quad \forall k$$

Let the FIR filter with transfer function  $W(z)$  in Figure 9.2 be characterized by the following difference equation or input-output relationship:

$$\hat{v}_1(n) = w(0)v_2(n) + w(1)v_2(n-1) + \cdots + w(m-1)v_2(n-m+1)$$

If we introduce the vector quantities:

$$\mathbf{w} = \begin{bmatrix} w(0) \\ w(1) \\ \vdots \\ w(m-1) \end{bmatrix} \quad \mathbf{v}_2(\mathbf{n}) = \begin{bmatrix} v_2(n) \\ v_2(n-1) \\ \vdots \\ v_2(n-m+1) \end{bmatrix}$$

the signal  $\hat{v}_1(n)$  can then compactly be written as:

$$\hat{v}_1(n) = \mathbf{w}^H \mathbf{v}_2(\mathbf{n}) \quad (9.19)$$

Based on the above notation and assumptions, the FIR Active Noise Cancellation Problem determines the FIR filter coefficients  $\mathbf{w}$  by minimizing the following criterion.

$$J(\mathbf{w}) = E[|\mathbf{w}^H \mathbf{x}(\mathbf{n}) - v_1(n)|^2] \quad (9.20)$$

The minimization of this criterion, which is again a minimum variance criterion on the error signal  $v_1(n) - \hat{v}_1(n)$ , can be done following the proof of Theorem 9.1. For this particular case and using the Auto-correlation function  $r_{v_2}(k)$  of  $v_2(n)$  and the Cross-correlation function  $r_{v_1 v_2}(k)$ , the Wiener-Hopf equation that determine the optimal filter coefficients  $\hat{\mathbf{w}}$  is given as:

$$\begin{bmatrix} r_{v_2}(0) & r_{v_2}^*(1) & \cdots & r_{v_2}^*(m-1) \\ r_{v_2}(1) & r_{v_2}(0) & & r_{v_2}^*(m-2) \\ \vdots & & \ddots & \\ r_{v_2}(m-1) & r_{v_2}(m-2) & \cdots & r_{v_2}(0) \end{bmatrix} \begin{bmatrix} \hat{w}(0) \\ \hat{w}(1) \\ \vdots \\ \hat{w}(m-1) \end{bmatrix} = \begin{bmatrix} r_{v_1 v_2}(0) \\ r_{v_1 v_2}(1) \\ \vdots \\ r_{v_1 v_2}(m-1) \end{bmatrix}$$

In a realistic scenario it may be assumed that the signals  $v_2(n)$  and  $y(n)$  can be measured. For that reason it would be difficult to have information on the

Cross-correlation function  $r_{v_1 v_2}(k)$ . However making use of the assumption that stochastic processes  $d(n)$  and  $v_2(n)$  are uncorrelated, we have that:

$$E[v_1(n)v_2^*(n-k)] = E[y(n)v_2^*(n-k)]$$

This is also a plausible assumption, as we want the stochastic process to be correlated to the additive noise  $v_1(n)$ , and the latter may also be assumed to be uncorrelated to the desired signal  $d(n)$ . Based on the equality of  $r_{v_1 v_2}(k)$  with  $r_{y v_2}(k)$  the Wiener-Hopf equation for solving the FIR Active Noise Cancellation Problem are:

$$\begin{bmatrix} r_{v_2}(0) & r_{v_2}^*(1) & \cdots & r_{v_2}^*(m-1) \\ r_{v_2}(1) & r_{v_2}(0) & & r_{v_2}^*(m-2) \\ \vdots & & \ddots & \\ r_{v_2}(m-1) & r_{v_2}(m-2) & \cdots & r_{v_2}(0) \end{bmatrix} \begin{bmatrix} \hat{w}(0) \\ \hat{w}(1) \\ \vdots \\ \hat{w}(m-1) \end{bmatrix} = \begin{bmatrix} r_{y v_2}(0) \\ r_{y v_2}(1) \\ \vdots \\ r_{y v_2}(m-1) \end{bmatrix} \quad (9.21)$$

And therefore the following values of its Auto-and Cross-correlation functions are necessary:

$$r_{v_2}(0), r_{v_2}(1), \dots, r_{v_2}(m-1), r_{y v_2}(0), r_{y v_2}(1), \dots, r_{y v_2}(m-1)$$

Assuming a (finite length) realization of the stochastic processes  $v_2(n)$  and  $y(n)$ , these coefficients may be estimated under the assumption of Auto- (and Cross-) correlation ergodicity. The Cross-correlation ergodicity is a minor extension of the notion of Auto-Correlation Ergodicity discussed in Sections 5.2.7 to the Cross-Correlation function between two stochastic processes.

## 9.4 Minimum Variance IIR Wiener Filter

### 9.4.1 The Generic Problem Formulation

For the IIR Wiener filter design problem we will consider the filter  $W(z)$  in the generic Block-scheme in Figure 9.1 to have two different forms. For the general IIR case, we consider it to have the following form:

$$W(z; w(k)) = \sum_{k=-\infty}^{\infty} w(k)z^{-k} \quad (9.22)$$

Where it is assumed that this z-transform has a ROC that contains the unit circle  $\Gamma$  in the complex plane. When  $W(z; w(k))$  is rational in the complex variable  $z$ , it would mean that the transfer function does not contain poles on the unit circle.

For the case where we assume *in addition* that the filter  $W(z; w(k))$  is causal, we consider  $W(z; w(k))$  to be restricted in this section as:

$$W_c(z; w(k)) = \sum_{k=0}^{\infty} w(k)z^{-k} \quad (9.23)$$

The parameters  $w(j)$  in both parametrized filters in (9.22) and (9.23) are to be designed to yield a minimum variance filter design, minimizing the variance of the error signal  $e(n; w(k))$  in Figure 9.1. This is again expressed a criterion, and results in the *Minimum Variance IIR Wiener Problem*

$$\min_{w(k)} J(w(k)) \quad \text{for } J(w(k)) = E[|e(n; w(k))|^2] \quad (9.24)$$

with for the case  $k \in \mathbb{Z}$ , i.e. when the filter  $W(z; w(k))$  is given as in (9.22), this will simply be referred to as the *mixed causal, anti-causal IIR Wiener filter* and for the case  $k \in \mathbb{Z}_+$ , i.e. when the filter  $W(z; w(k))$  is given as in (9.23), this will simply be referred to as the *causal IIR Wiener filter*.

#### 9.4.2 The mixed causal, anti-causal IIR Wiener Filter

The solution to the mixed causal, anti-causal IIR Wiener Problem is summarized in the following Theorem.

**Theorem 9.2** (The mixed causal, anti-causal IIR Wiener Filter). *Let the conditions as stipulated in the generic problem formulation given in Section 9.4.1 hold for the filter  $W(z)$  to be defined as in (9.22), let the stochastic processes  $x(n)$  and  $d(n)$  be jointly zero mean and WSS, with the following Power and cross-spectra:*

$$P_x(e^{j\omega}) > 0 \quad P_{dx}(e^{j\omega}) = \sum_{k=-\infty}^{\infty} r_{dx}(k) e^{-j\omega k} \quad (9.25)$$

then the solution to the mixed causal, anti-causal IIR Wiener Problem (9.24) is given in terms of the optimal filter coefficient vector  $\hat{W}(z)$  and value of the criterion at this optimum as:

$$\hat{W}(z) = P_{dx}(z) P_x(z)^{-1} \quad (\text{The Wiener-Hopf equations}) \quad (9.26)$$

$$J(\hat{W}(z)) = \frac{1}{2\pi} \int_{-\pi}^{\pi} P_d(e^{j\omega}) - P_{dx}(e^{j\omega}) P_x(e^{j\omega})^{-1} P_{dx}^*(e^{j\omega}) d\omega \quad (9.27)$$

*Proof.* Following the proof of Theorem 9.1 we arrive in the same way at the following *Orthogonality Condition*:

$$E[e(n; \hat{w}(k)) x^*(n - k)] = 0 \quad -\infty < k < \infty \quad (9.28)$$

Using the expression for the error  $e(n; w(k))$  given as:

$$e(n; w(k)) = d(n) - \sum_{\ell=-\infty}^{\infty} w(k) x(n - \ell) \quad -\infty < k < \infty \quad (9.29)$$

and using the linearity of the Expectation Operator  $E[\cdot]$ , the Orthogonality condition becomes:

$$E[d(n) x^*(n - k)] - E\left[ \sum_{\ell=-\infty}^{\infty} \hat{w}(\ell) x(n - \ell) x^*(n - k) \right] = 0 \quad -\infty < k < \infty$$

With the definition of the cross-correlation function  $r_{dx}(k) = E[d(n)x^*(n-k)]$  and the Auto-correlation function  $r_x(k) = E[x(n)x^*(n-k)]$ , we can write this equation as:

$$r_{dx}(k) - \sum_{\ell=-\infty}^{\infty} \hat{w}(\ell)r_x(k-\ell) = 0 \quad -\infty < k < \infty$$

Taken the DTFT yields:

$$P_{dx}(e^{j\omega}) - \hat{W}(e^{j\omega})P_x(e^{j\omega}) = 0$$

And since  $P_x(e^{j\omega}) > 0$  and considering the z-transform, the result for  $\hat{W}(z)$  in (9.26) results.

To find the optimal value of the criterion  $J(\hat{W}(z))$ , we again make use of the Orthogonality condition (9.28) and the expression of  $e(n; w(j))$  given in (9.29) to express this optimum as,

$$J(\hat{W}(z)) = E[e(n; \hat{w}(k))d^*(n)]$$

Again using the expression for  $e(n; w(j))$  in (9.29), this can be written as,

$$J(\hat{W}(z)) = E\left[\left(d(n) - \sum_{\ell=-\infty}^{\infty} \hat{w}(\ell)x(n-\ell)\right)d^*(n)\right]$$

or,

$$J(\hat{W}(z)) = r_d(0) - \sum_{\ell=-\infty}^{\infty} \hat{w}(\ell) \underbrace{E[d(n)x^*(n-\ell)]^*}_{r_{dx}^*(\ell)} \quad (9.30)$$

If we now define the signal  $\alpha(k)$  as  $\sum_{\ell=-\infty}^{\infty} \hat{w}(\ell)r_{dx}^*(k+\ell)$ , then we have the following relationships between their DTFTs,

$$\mathcal{F}(\alpha(n)) = \hat{W}(e^{j\omega})P_{dx}^*(e^{j\omega})$$

Using the result of the inverse DTFT given by (2.10) yields,

$$\alpha(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \hat{W}(e^{j\omega})P_{dx}^*(e^{j\omega})d\omega$$

and since  $\alpha(0)$  equals  $\sum_{\ell=-\infty}^{\infty} \hat{w}(\ell)r_{dx}^*(\ell)$ , and by using (5.11) (which is again using the definition of the inverse DTFT), (9.30) can be written as (9.27).  $\square$

### 9.4.3 Application of the mixed causal, anti-causal IIR Wiener filter to Denoising

We consider the same signal configuration as in Figure 9.1 and assumptions on the signals as in subsection 9.3.3. However now we investigate an IIR filter design. The denoising problem fits in the generic problem formulation for the

stochastic processes  $x(n)$  and  $d(n)$  as related in (9.15) with the noise  $v(n)$  satisfying (9.14). For the solution given in Theorem 9.2 we need the the following info of these stochastic processes,

$$P_x(z) \quad P_{dx}(z)$$

From the corresponding Auto- and Cross correlation functions given in (9.16) and (9.17) resp. these spectra can be written as

$$P_x(z) = P_d(z) + P_v(z) \quad \text{and} \quad P_{dx}(z) = P_d(z)$$

Therefore the minimum variance mixed causal, anti-causal IIR Wiener filter for the denoising problem equals,

$$\hat{W}(z) = \frac{P_d(z)}{P_d(z) + P_v(z)} \quad (9.31)$$

With the optimal filter written as in (9.26), the optimal value of the cost  $J(\hat{W}(z))$  in (9.27) can also be written as,

$$J(\hat{W}(z)) = \frac{1}{2\pi} \int_{-\pi}^{\pi} P_d(e^{j\omega}) - \hat{W}(e^{j\omega}) P_{dx}^*(e^{j\omega}) d\omega$$

For the specific optimal filter for the denoising problem, this cost becomes,

$$J(\hat{W}(z)) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \left[ P_d(e^{j\omega}) - \frac{P_d(e^{j\omega})}{P_d(e^{j\omega}) + P_v(e^{j\omega})} P_d^*(e^{j\omega}) \right] d\omega$$

or,

$$\begin{aligned} J(\hat{W}(z)) &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left[ \frac{P_d(e^{j\omega})^2 + P_d(e^{j\omega})P_v(e^{j\omega}) - P_d(e^{j\omega})^2}{P_d(e^{j\omega}) + P_v(e^{j\omega})} \right] d\omega \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left[ P_v(e^{j\omega}) \frac{P_d(e^{j\omega})}{P_d(e^{j\omega}) + P_v(e^{j\omega})} \right] d\omega \end{aligned} \quad (9.32)$$

The last expression shows that the optimal cost  $J(\hat{W}(z))$  becomes zero when the power spectra  $P_v(e^{j\omega})$  (of the noise) and  $P_d(e^{j\omega})$  (of the desired signal) *do not overlap*. As the cost represents the variance of the error between the desired signal and the reconstructed one, it means that under this circumstance a *perfect denoising* is possible.

#### 9.4.4 Application of the mixed causal, anti-causal IIR Wiener filter to Deconvolution

The deconvolution problem is defined for the stochastic processes  $x(n)$  and  $d(n)$  of the general problem formulation in subsection 9.4.1 given as:

$$x(n) = \sum_{\ell=-\infty}^{\infty} g(\ell)d(n-\ell) + r(n) \quad (9.33)$$



with the additive noise  $r(n)$  uncorrelated from the desired signal  $d(n)$ , given as

$$E[d(n)r^*(n-k)] = 0 \quad \forall k \quad (9.34)$$

For this situation the calculation of the minimum variance mixed causal, anti-causal IIR Wiener requires the spectra,

$$P_{dx}(z) \quad P_x(z)$$

for the signal model given as in (9.33) and the assumption (9.34), these spectra become,

$$\begin{aligned} P_{dx}(z) &= G^*(1/z^*)P_d(z) \\ P_x(z) &= G(z)G^*(1/z^*)P_d(z) + P_r(z) \end{aligned}$$

And the optimal filter  $\hat{W}(z)$  is given as:

$$\hat{W}(z) = \frac{G^*(1/z^*)P_d(z)}{G(z)G^*(1/z^*)P_d(z) + P_r(z)} \quad (9.35)$$

#### 9.4.5 The Causal IIR Wiener Filter

##### A possible breakdown in a straightforward application of Theorem 9.2 to the causal case

When restricting the Wiener filter to being causal as in (9.23) many steps for the mixed causal, anti-causal optimal Wiener design can be repeated. However at some point this analogy breaks down. In this subsection we show where this happens and that a different approach is needed.

Similar to the first four equations in the proof of Theorem 9.2 we can arrive at the following related equations for the Minimum Variance causal Wiener filter design. The orthogonality condition now reads:

$$E[e(n; \hat{w}(k))x^*(n-k)] = 0 \quad 0 \leq k < \infty \quad (9.36)$$

Using the expression for the error  $e(n; w(k))$ , now given as:

$$e(n; w(k)) = d(n) - \sum_{\ell=0}^{\infty} w(k)x(n-\ell) \quad (9.37)$$

and the linearity of the Expectation Operator  $E[\cdot]$ , the Orthogonality condition becomes:

$$E[d(n)x^*(n-k)] - E\left[\sum_{\ell=0}^{\infty} \hat{w}(\ell)x(n-\ell)x^*(n-k)\right] = 0 \quad 0 \leq k < \infty$$

With the definition of the cross-correlation function  $r_{dx}(k) = E[d(n)x^*(n-k)]$  and the Auto-correlation function  $r_x(k) = E[x(n)x^*(n-k)]$ , we can write this equation as:

$$r_{dx}(k) - \sum_{\ell=0}^{\infty} \hat{w}(\ell)r_x(k-\ell) = 0 \quad 0 \leq k < \infty \quad (9.38)$$

This equation is only defined for  $0 \leq k < \infty$ . However, if we would pretend it to hold for  $-\infty < k < 0$  as well, we could take the z-transform and subsequently obtain,

$$\begin{aligned} P_{dx}(z) - \sum_{k=-\infty}^{\infty} \sum_{\ell=0}^{\infty} \hat{w}(\ell) r_x(k-\ell) z^{-k} &= 0 \\ P_{dx}(z) - \sum_{\ell=0}^{\infty} \hat{w}(\ell) \sum_{k=-\infty}^{\infty} r_x(k-\ell) z^{-k} &= 0 \\ P_{dx}(z) - \left( \sum_{\ell=0}^{\infty} \hat{w}(\ell) z^{-\ell} \right) P_x(z) &= 0 \end{aligned}$$

However this breaks down since equality can not hold (in general) since  $P_{dx}(z)P_x(z)^{-1}$  is mixed causal, anti-causal. And as such can never be matched with just the causal term  $\left( \sum_{\ell=0}^{\infty} \hat{w}(\ell) z^{-\ell} \right)$ .

For that reason a different approach is needed. This will be done in two steps, following the outline in [1]. First we try to solve (9.38) when the stochastic process  $x(n)$  is ZMWN. Though it is a very special case of little practical value, it paves the way to the full general solution of the Minimum Variance IIR causal Wiener filter design problem. These two steps are discussed next.

### Solution for $x(n)$ ZMWN

For the case  $x(n)$  is ZMWN with variance  $\sigma_x^2$ , its Auto-correlation function is  $r_x(k) = \sigma_x^2 \Delta(k)$ . Therefore (9.38) simply becomes,

$$r_{dx}(k) = \sigma_x^2 \hat{w}(k) \quad 0 \leq k < \infty$$

Again using the notation to indicate the causal part of a z-transform on page 25, the z-transform of the optimal Wiener filter in this case is,

$$\hat{W}(z) = \frac{[P_{dx}]_+}{\sigma_x^2} \quad (9.39)$$

### Solution to the general Minimum Variance IIR causal Wiener filter design problem

Following the approach outlined in [1], we can use the result of the previous subsection by whitening the stochastic process  $x(n)$ . Hereby we make use of the *whitening filter* based on the spectral factorization as discussed in Section 7.3.3. Let the Spectral factorization of  $x(n)$  be given as:

$$P_x(z) = \sigma_0^2 Q(z) Q^*(1/z^*)$$

with  $Q(z)$  minimum phase. Then filtering  $x(n)$  by the filter with transfer function  $\frac{1}{\sigma_0 Q(z)}$  yields as depicted in the left part of Figure 9.3 the white noise signal

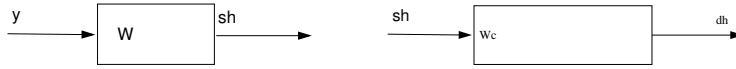
$\epsilon(n)$ . Based on the previous subsection we can state that the solution to the following Causal IIR Wiener filter design problem:

$$\min_{w_c(\ell)} E[|d(n) - \sum_{\ell=0}^{\infty} w_c(\ell)\epsilon(n-\ell)|^2] \quad (9.40)$$

is given by,

$$\hat{W}_c^w = [P_{d\epsilon}(z)]_+ \quad (9.41)$$

This filter provides a minimum variance estimate of  $d(n)$  denoted as  $\hat{d}(n; w(k))$  in the right part of Figure 9.3. To represent the filter  $\hat{W}_c$  given in (9.41) in terms



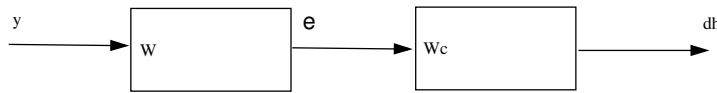
**Figure 9.3:** The whitening of a generic stochastic process  $x(n)$  (Left) and the Minimum Variance causal Wiener filter from the whitened stochastic process  $\epsilon(n)$  (Right).

of quantities (such as Power and/or Cross spectra) that can be derived from the original stochastic process  $d(n)$  and  $x(n)$ , we make use of the solution to Exercise 6.5. Consider the relationship between  $\epsilon(n)$  and  $x(n)$  as depicted in Figure 9.3 (Left) then the Cross-spectrum  $P_{d\epsilon}(z)$  is related to  $P_{dx}(z)$  by,

$$P_{d\epsilon}(z) = \frac{1}{\sigma_0^2 Q^*(1/z^*)} P_{dx}(z)$$

Combining the whitening filtering with the minimum variance IIR Causal Wiener filter using this whitened signal as given in (9.41), and depicted in Figure 9.4, yields the following cascaded filter that solves minimum variance IIR Causal Wiener filter design problem,

$$\hat{W}_c(z) = \frac{1}{\sigma_0^2 Q(z)} \left[ \frac{P_{dx}(z)}{Q^*(1/z^*)} \right]_+ \quad (9.42)$$



**Figure 9.4:** The cascade of the whitening filter and the Minimum Variance causal Wiener filter from the whitened stochastic process  $\epsilon(n)$  of Figure 9.3.

## 9.5 Example

In this section we review the theory outlined in Section 9.4 by applying the theory to a concrete example. This example considers the denoising problem for the signal  $d(n)$  given by the AR model of order 1:

$$d(n) - 0.8d(n-1) = 0.6w(n) \quad (9.43)$$

with  $w(n)$  a zero-mean WSS white noise with unit variance. A disturbed signal  $x(n)$  can be recorded and this signal is related to  $d(n)$  in the following way:

$$x(n) = d(n) + v(n) \quad (9.44)$$

with  $v(n)$  a zero-mean WSS white noise with unit variance that is statistically uncorrelated from  $w(\ell)$  for all  $n, \ell$ . The additive noise  $v(n)$  also satisfies,

$$E[d(n)v(n-k)] = 0 \quad \forall k \quad (9.45)$$

### 9.5.1 Calculating the mixed causal, anti-causal IIR Wiener filter

Then the mixed causal, anti-causal IIR Wiener filter (9.31) derived in Section 9.4.3 requires the evaluation of  $P_d(z)$  and  $P_v(z)$ . For the above given signal information these are:

$$\begin{aligned} P_d(z) &= \left| \frac{0.6\sigma_w^2}{1-0.8z^{-1}} \right|^2 = \frac{0.36}{(1-0.8z^{-1})(1-0.8z)} \\ P_v(z) &= \sigma_v^2 = 1 \end{aligned} \quad (9.46)$$

Substituting these results into (9.31) yields for this example,

$$\begin{aligned} \hat{W}(z) &= \frac{0.36}{0.36 + (1-0.8z^{-1})(1-0.8z)} \\ &= \frac{0.36}{2 - 0.8z^{-1} - 0.8z} \\ &= \frac{0.36}{1.6(1-0.5z)(1-0.5z^{-1})} \\ &= \frac{0.225}{(1-0.5z)(1-0.5z^{-1})} \\ &= \frac{0.3(1-\alpha^2)}{(1-\alpha z)(1-\alpha z^{-1})} \quad \alpha = 0.5 \end{aligned} \quad (9.47)$$

To compute the cost given by (9.32) we make use of Theorem 2.2 for  $X(z) = \hat{W}(z)$ . Then the cost becomes,

$$\begin{aligned} J(\hat{W}(z)) &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \hat{W}(e^{j\omega}) d\omega \\ &= \frac{1}{2\pi j} \oint_{\Gamma} \frac{\hat{W}(z)}{z} dz \\ &= \hat{w}(0) \end{aligned} \quad (9.48)$$

With  $\hat{w}(n)$  the (two-sided) impulse response of the filter with transfer function  $\hat{W}(z)$ . Using Table 2.4 we find that,

$$\hat{w}(n) = 0.3 \frac{1}{2}^{|n|}$$

And therefore the cost is given by,

$J(\hat{W}(z)) = 0.3$

### 9.5.2 Calculating the causal IIR Wiener filter

For computing the causal Wiener filter for the signals defined at the beginning of Section 9.5, we need to compute the spectral factorization of  $P_x(z)$ . With the latter transfer function given as,

$$\begin{aligned} P_x(z) &= P_d(z) + P_v(z) \\ &= \frac{0.36}{(1 - 0.8z^{-1})(1 - 0.8z)} + 1 \\ &= \frac{2 - 0.8z^{-1} - 0.8z}{(1 - 0.8z^{-1})(1 - 0.8z)} \end{aligned} \quad (9.49)$$

Using the results in (9.47), we can express  $P_x(z)$  also as,

$$P_x(z) = 1.6 \frac{(1 - 0.5z^{-1})(1 - 0.5z)}{(1 - 0.8z^{-1})(1 - 0.8z)}$$

Which displays the spectral factorization directly. So,

$$\sigma_0^2 = 1.6 \quad Q(z) = \frac{(1 - 0.5z^{-1})}{(1 - 0.8z^{-1})}$$

As a consequence the causal Wiener filter given by (9.42) becomes,

$$\hat{W}_c(z) = \frac{(1 - 0.8z^{-1})}{1.6(1 - 0.5z^{-1})} \left[ \frac{0.36}{(1 - 0.8z^{-1})(1 - 0.8z)} \frac{(1 - 0.8z)}{(1 - 0.5z)} \right]_+ \quad (9.50)$$

To compute the causal part of  $\frac{0.36}{(1 - 0.8z^{-1})(1 - 0.5z)}$  we perform its partial fraction expansion. The latter is outlined in the Appendix of [2]. The result is,

$$\frac{0.36}{(1 - 0.8z^{-1})(1 - 0.5z)} = \frac{0.6}{1 - 0.8z^{-1}} + \frac{0.3z}{1 - 0.5z}$$

Therefore,

$$\begin{aligned} \hat{W}_c(z) &= \frac{(1 - 0.8z^{-1})}{1.6(1 - 0.5z^{-1})} \frac{0.6}{(1 - 0.8z^{-1})} \\ &= \frac{0.375}{1 - 0.5z^{-1}} \end{aligned} \quad (9.51)$$

The value of the cost function is given by,

$$\begin{aligned} J(\hat{W}_c(z)) &= E\left[\left(d(n) - \sum_{\ell=0}^{\infty} \hat{w}_c(\ell)x(n - \ell)\right)d(n)\right] \\ &= r_d(0) - \sum_{\ell=0}^{\infty} \hat{w}_c(\ell)r_d(\ell) \end{aligned} \quad (9.52)$$

With  $w_c(\ell) = 0.375 \cdot 0.5^\ell$  and  $r_d(\ell) = 0.8^\ell$  (for  $\ell \geq 0$ ), we obtain,

$$\boxed{J(\hat{W}_c(z)) = 1 - 0.375 \sum_{\ell=0}^{\infty} 0.5^\ell 0.8^\ell = 0.375} \quad (9.53)$$

---

## References

- [1] M. H. Hayes, *Statistical Digital Signal Processing and Modeling*. New York: John Wiley and Sons, 1996.
  - [2] A. V. Oppenheim and A. S. Willsky, *Signals and Systems*. Upper Saddle River, New Jersey: Prentice-Hall, second ed., 1997.
- 

## Exercises

**Exercise 9.1** Consider the desired zero-mean WSS stochastic process  $d(n)$ . The ZMWN noise  $v(n)$  that is uncorrelated with  $d(n)$  and satisfies the condition (9.14) is an additive perturbation to the measurement of  $d(n)$  as:

$$x(n) = d(n) + v(n)$$

The auto-correlation of  $d(n)$  is given and denoted as  $r_d(k)$  and the variance  $v(n)$  denoted as  $\sigma_v^2$  is also given. Then we seek to estimate  $d(n+1)$  (a one-step ahead prediction) by the following FIR model of order 2:

$$\hat{d}(n+1; \mathbf{w}) = w(0)x(n) + w(1)x(n-1)$$

such that the following variance is minimized.

$$J(\mathbf{w}) = E[|d(n+1) - \hat{d}(n+1; \mathbf{w})|^2]$$

with  $\mathbf{w} = \begin{bmatrix} w(0) \\ w(1) \end{bmatrix}$ .

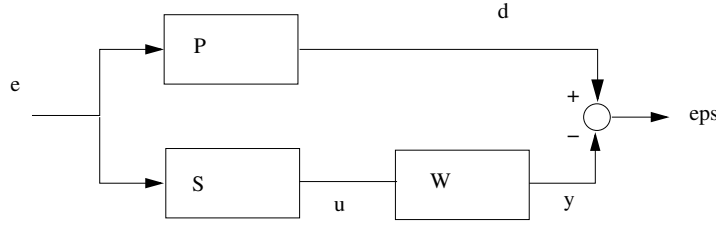
The solution to this problem needs to address the following subproblems:

- (a) Determine (as a function of the given information) the Wiener Hopf equations for this Minimum Variance FIR Wiener filter problem when  $r_d(k) = 0.8^{|k|}$  and  $\sigma_v^2 = 1$ .
- (b) Determine the optimal transfer function  $W(z)$  from the optimal filter coefficients derived from the Wiener Hopf equations derived in Part (a) of this exercise.
- (c) Determine the optimal value of the criterium  $J(\hat{\mathbf{w}})$ .
- (d) What is the value of the optimal filter coefficients if  $\sigma_v^2$  becomes zero.

**Exercise 9.2** Consider the noise cancellation problem depicted in Figure 9.5. The transfer functions  $P(z)$  and  $S(z)$  in Figure 9.5 are given by:

$$P(z) = 1 - 1.2z^{-1} \quad S(z) = 1 - 2.4z^{-1} + 1.35z^{-2}$$

The stochastic process  $e(k)$  is zero-mean white noise with unknown variance  $\sigma_e^2$ . The goal is to design a filter  $W(z)$  such that the variance of  $\epsilon(k) = d(k) - y(k)$  is minimized.



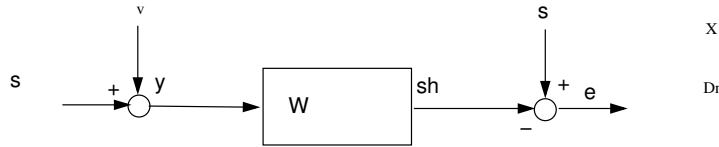
**Figure 9.5:** A blockscheme of the Noise Cancellation problem considered in Question 9.2.

- Determine the IIR-noncausal optimal Wiener filter solution for  $W(z)$  and call this optimal filter  $W_1(z)$ .
- Determine the resulting variance of the residual  $\epsilon(k)$  by use of this IIR-noncausal Wiener filter  $W_1(z)$ .
- Give a stable difference equation (i.e. a recursion in time) to execute the IIR-noncausal Wiener filter  $W_1(z)$  on the sequence  $u(k)$  when this sequence is available from  $k = -\infty, \dots, \infty$ .

**Exercise 9.3** Consider the denoising problem depicted in Figure 9.6. Here the real, Wide Sense Stationary (WSS) Random Process (RP)  $d(n)$  satisfies,

$$d(n) = 0.8d(n-1) + r(n)$$

where  $r(n)$  is ZMWN(0.36). The additive noise  $v(n)$  that disturbs the measurement  $x(n) = d(n) + v(n)$  in Figure 9.6 is ZMWN(1) and is uncorrelated from  $r(\ell)$  for all  $n, \ell$ . A causal filter  $W(z)$  has been designed with one free



**Figure 9.6:** The Denoising problem in Exercise 9.3

parameter  $a$  with the following transfer function:

$$W(z) = \frac{a}{1 - 0.5z^{-1}} \quad (9.54)$$

The goal is to determine the parameter  $a$  such that the following criterium is minimized

$$\min_a E[|d(n) - \hat{d}(n)|^2]$$

For that purpose the following questions need to be addressed to solve this exercise.

- Determine the Autocorrelaton function  $r_d(k)$  of the RP  $d(n)$ .
- Give an analytic expression for the cross-correlation function  $r_{\hat{d},d}(k)$  and determine its value for  $k = 0$  as a function of the parameter  $a$ .

- (c) Determine the variance  $E[\hat{d}(n)^2]$  as a function of the parameter  $a$ .
- (d) Determine the numerical value of the parameter  $a$  that minimizes the criterium

$$E[|d(n) - \hat{d}(n)|^2].$$

**Exercise 9.4** In a Compact Disc player, the position of the laserspot on the disc is controlled by measuring the displacement of the spot position compared to the desired position on the right track. This displacement can be modelled as a discrete-time stationary stochastic process  $x$ , described by

$$x(t) = 0.5x(t-1) + e(t)$$

where  $e$  is a white noise process with variance 1.

Suppose the displacement  $x(t)$  can only be measured with a disturbance  $v(t)$ :

$$y(t) = x(t) + v(t)$$

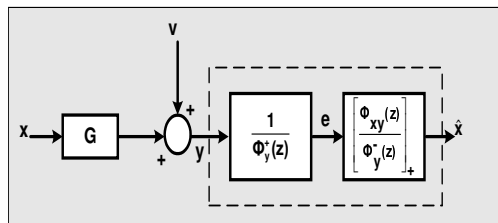
where  $v$  is a white noise process with variance  $8/7$ , and  $x$  and  $v$  are uncorrelated.

- (a) Calculate the power spectrum  $P_x(\omega)$  of the process  $x$ .
- (b) Suppose we only have one measurement  $y(t)$ . The linear mean square error estimator of  $x(t)$  based on this measurement can then be denoted as:

$$\hat{x}(t) = a_{ms} + b_{ms}y(t)$$

Compute the coefficients  $a_{ms}$  and  $b_{ms}$  of this estimator.

- (c) Compute the transfer function (in the  $Z$ -domain) of the optimal, linear, infinite impulse response non-causal filter estimating  $x(t)$  based on the measurements  $y(n), n \in (-\infty, \infty)$  (the non-causal Wiener filter).
- (d) Compute the transfer function (in the  $Z$ -domain) of the optimal, linear, infinite impulse response causal filter estimating  $x(t)$  based on the measurements  $y(n), n \in (-\infty, \infty)$  (the causal Wiener filter). The causal Wiener filter is depicted schematically in Figure 9.7.



**Figure 9.7:** Causal Wiener filter



**Exercise 9.5** Consider a process  $x$  that can be modeled as a stationary stochastic process with expected value 0 and auto-correlation function:

$$r_x(\tau) = 0.5^{|\tau|}.$$

Furthermore, assume that  $x(t)$  can only be measured under the influence of a disturbance  $v(t)$ :

$$y(t) = x(t) + v(t)$$

where  $v$  is a white noise process with variance 0.5, while  $x$  and  $v$  are uncorrelated.

- (a) Determine the transfer function ( $z$ -domain) of the IIR non-causal Wiener filter that provides an estimate of  $x(t)$  based on the observations  $y(n)$ ,  $n \in (-\infty, \infty)$ .
- (b) Determine the coefficients  $a$  and  $b$  of the transfer function

$$H_c(z) = \frac{a}{1 - bz^{-1}}$$

of the IIR causal Wiener filter that estimated  $x(t)$  based on the observations  $y(n)$ ,  $n \in (-\infty, t]$ .



# Subject Index

- Absolute Summable, 19
- Anti-Causal LTI Systems, 25
- Anti-Causal Part of Transfer Function, 25
- AR model, 112
- ARMA model, 111
- Auto-correlation (function), 86
- Auto-correlation Ergodic, 94
- Auto-correlation Matrix, 91
- Auto-covariance (function), 86
  
- Bernoulli process, 99
- bias, 49
- BIBO Stability, 24
- BLUE, 63
  
- Cauchy Schwartz inequality, 42
- Cauchy's Integral Formula, 22
- Causal LTI Systems, 25
- Causal Part of Transfer Function, 25
- Conjugate Symmetric, 20
- consistency, 50
- Continuous Random Variable, 36
- Convolution, 20
- Coprime, 128
- Correlation, 41
- Correlation Coefficient, 41
- Covariance, 41
- Covariance matrix LLSQ problem, 159
- Cramér-Rao bound, 64
- Cross Spectrum, 97
- Cross-correlation (function), 87
- Cross-covariance (function), 87
- Cybernetics, 167
  
- Data Equation, 153
- Dependent Variable, 154
- dependent variables, 53
- Discrete Random Variable, 36
- Discrete-Time Fourier Transform (DTFT), 19
- Discrete-Time Signals, 18
- Discrete-Time Systems, 22
  
- efficiency, 50
- Embedding of finite time series in an infinite one, 137
- Ensemble Average, 38
- Ergodic in the mean, 92
- Ergodicity, 92
  
- estimator, 48
  - asymptotically unbiased, 49
  - efficient, 50
  - MSE, 49
  - unbiased, 49
  - variance, 49
- Example causal IIR Wiener Filter, 183
- Example mixed causal, anti-causal IIR Wiener Filter, 182
- Existence z-transform, 21
- Expectation Operator, 38
- Expected Value, 38
  
- Finite Impulse Response (FIR), 166
- FIR Active Noise Cancellation Problem, 174
- First Moment, 38
- Fisher Information Matrix, 64
- Forward Modeling, 104, 126
  
- Gaussian Processes, 88
- Gaussian Random Variable(s), 43
- Geometric Interpretation, 157
  
- Hermitian Transpose, 155
  
- IIR Causal Wiener Filter, 181
- Independent Random Variables, 42
- Independent Variables, 154
- independent variables, 53
- Infinite Impulse Response (IIR), 166
- Inverse DTFT, 20
- Inverse of an LTI system, 27
  
- Joint Density Function, 40
- Joint Distribution Function, 40
  
- Lag, 89
- least-squares, 53
  - weighted, 61
- likelihood function, 69
- Linear Least Squares for AR models, 150
- Linear Operator, 41
- Linear Discrete Time Systems, 23
- linear regresion
  - least squares solution, 54, 56
- Linear Regression, 152
- linear regression, 52

- MA model, 112
- Markov estimator, 63
- maximum likelihood estimator, 69
  - properties, 74
- Mean, 38
- Mean of a Stochastic Process, 86
- Mean Square Error, 40
- mean squared error (MSE)
  - of estimator, 49
- Means Square Value, 39
- Minimum Variance Estimation, 150
- minimum variance estimator, 64
- Minimum Variance Filter design, 167
- Minimum Variance FIR Wiener Problem, 169
- Minimum-Phase, 27
- Modified Yule-Walker equation, 139
- MSE, 49
  
- Nonlinear Transformation RV, 40
- Norbert Wiener, 167
- Normal Equations, 156
- normal equations, 54
  
- One-step Ahead Prediction, 173
- Order, 169
- Orthogonal Random Variables, 42
- Orthogonal Stochastic Processes, 87
- Orthogonality Condition, 157
  
- Partial Fraction Expansion, 26
- Perfect Denoising, 178
- Periodogram, 97
- Positive Real, 129
- Power Spectrum, 95
- Power Spectrum AR, 114
- Power Spectrum ARMA, 113
- Power Spectrum MA, 114
- Probability Distribution Function (PDF), 37
- Probability Density Function (pdf), 37
- Properties z-transform, 21
  
- Quadratic Optimization Problems, 152
  
- Rational, 128
- Realization of a stochastic Process, 85
- Recipe for Spectral Factorization, 131
- Region of Convergence (ROC), 21
- regressor, 53
- Regular Stochastic Processes, 96
  
- Second Moment, 39
- Shaping Filter, 126
- Simulation, 126
- Spectral Factor, 129
- statistic, 48
- Stochastic Processes, 84
- Synthesis problems, 166
- System Identification, 126
  
- The inverse AR-problem, 135
- The inverse ARMA-problem, 135
- The inverse MA-problem, 135
- Time-Invariant Discrete Time Systems, 23
  
- Unbiased Estimate, 159
- unbiased estimator, 49
- Uncorrelated Random Variables, 42
- Uncorrelated Stochastic Processes, 87
  
- Variance, 40
- variance
  - of estimator, 49
- Variance Complex Random Variable, 40
- Variance of a Stochastic Process, 86
  
- weighted least-squares, 61
- Whitening Filter, 133
- Wide Sense Stationarity, 89
- Wiener Filter, 166
  
- Yule-Walker equation, 115
  
- z-Transform, 21
- Zero-Mean Assumption, 87
- Zero-mean White Noise (ZMWN), 108