# SELECTED TOPICS IN PROBABILITY

Lecture notes by

Marco C. Campi Department of Information Engineering University of Brescia Italy

marco.campi@ing.unibs.it

January 7, 2021

**Foreword:** This text contains a treatment of some specific topics in probability theory. The selection of the topics has followed the goal of presenting in a self-contained manner all material that is necessary to understand the fundamental results in the theory of stochastic linear systems in discrete time, including prediction, filtering and control. In preparing these notes, I have put a significant amount of effort into combining requirements of conciseness and mathematical rigor with those of readability and clarity of presentation. I shall be grateful to anyone who will provide comments and suggestions on how to improve these notes.

# Contents

1	MEASURE SPACES AND INTEGRATION			
	1.1	Measurable spaces and measurable functions	3	
	1.2	Measures and measure spaces	6	
	1.3	Integration	9	
2	RANDOM VARIABLES			
	2.1	Random variables	15	
	2.2	Independence and uncorrelation	23	
	2.3	Characteristic functions	25	
	2.4	Gaussian random variables	31	
	2.5	Computing the density induced by a function	33	
3	STOCHASTIC CONVERGENCE			
	3.1	Probabilistic notions of convergence	37	
	3.2	Measurability of the limit of random variables	40	
	3.3	Limit under the sign of expectation	43	
	3.4	Convergence results for independent random variables	43	
	3.5	Weak convergence on $\mathbb{R}$	55	
4	THE PROJECTION THEOREM			
	4.1	Hilbert spaces	63	
	4.2	The projection theorem	74	
	4.3	Applications of the projection theorem	76	

# CONTENTS

5	COI	NDITIONAL EXPECTATION AND CONDITIONAL DENSITY	79
	5.1	Conditional expectation	79
	5.2	Conditional density	86
6	WII	DE-SENSE STATIONARY PROCESSES	91
	6.1	Definitions and examples	91
	6.2	Elementary spectral theory of stationary processes	95
	6.3	Spectral theory of stationary processes	95
	6.4	Multivariable stationary processes	105

1

# Chapter 1

# MEASURE SPACES AND INTEGRATION

# **1.1** Measurable spaces and measurable functions

#### **Measurable spaces**

The notion of measurable space makes use of the concept of  $\sigma$ -algebra. A  $\sigma$ -algebra is a collection of subsets of a given set, with certain set-theoretic properties specified in the following definition.

**DEFINITION 1.1** ( $\sigma$ -algebra) Given a set X, a collection  $\mathscr{X}$  of subsets of X is called a  $\sigma$ -algebra if

(a) 
$$X \in \mathscr{X}$$
;  
(b) if  $A_k \in \mathscr{X}$ ,  $k = 1, 2, ..., then \cup_{k=1}^{\infty} A_k \in \mathscr{X}$ ;  
(c) if  $A \in \mathscr{X}$ , then  $A^c \in \mathscr{X}$  ( $A^c$  is the complement of set A).  $\Box$ 

Condition (b) states that a  $\sigma$ -algebra is closed under countably infinite union, that is union of an infinite number of subsets  $A_k$ , where k runs over the integers. Since the empty set  $\emptyset$  equals  $X^c$  and  $X \in \mathscr{X}$  by (a), condition (c) implies that  $\emptyset \in \mathscr{X}$ . Taking  $A_{n+1} = A_{n+2} = \cdots = \emptyset$  in (b), we see that  $\bigcup_{k=1}^n A_k \in \mathscr{X}$ , that is  $\mathscr{X}$  is also closed under finite union. Moreover, a  $\sigma$ -algebra is also closed under intersection (finite or countably infinite), as it follows from relation  $\bigcap_k A_k = (\bigcup_k A_k^c)^c$ .

Given any collection  $\mathscr{A}$  of subsets of a set *X*, consider all  $\sigma$ -algebras containing  $\mathscr{A}$ . Their intersection (that is the collection of the sets that belong to all  $\sigma$ -algebras) is easily seen to be a  $\sigma$ -algebra too. It is called the minimal  $\sigma$ -algebra containing  $\mathscr{A}$  and is denoted by  $\sigma(\mathscr{A})$ . We can now introduce the notion of measurable space.

**DEFINITION 1.2 (measurable space)** A couple  $(X, \mathscr{X})$  where X is any set and  $\mathscr{X}$  is a  $\sigma$ -algebra of subsets of X is called a measurable space.

It is sometimes important to consider measurable spaces that are the product of other measurable spaces. This is formalized in the following definition.

**DEFINITION 1.3 (product measurable space)** Given the measurable spaces  $(X_k, \mathscr{X}_k), k = 1, 2, ..., n$ , their product measurable space is  $(X_1 \times X_2 \times \cdots \times X_n, \mathscr{X}_1 \otimes \mathscr{X}_2 \otimes \cdots \otimes \mathscr{X}_n)$ , where  $X_1 \times X_2 \times \cdots \times X_n$  is the direct product of the  $X_k$ 's, i.e. the set of ordered n-tuples  $(x_1, x_2, \cdots, x_n)$  with  $x_k \in X_k$ , k = 1, 2, ..., n, and  $\mathscr{X}_1 \otimes \mathscr{X}_2 \otimes \cdots \otimes \mathscr{X}_n$  is the direct product of the  $\mathscr{X}_k$ 's, i.e. the smallest  $\sigma$ -algebra in  $X_1 \times X_2 \times \cdots \times X_n$  that contains all sets of the form  $A_1 \times A_2 \times \cdots \times A_n = \{(x_1, x_2, \cdots, x_n) \text{ such that } x_k \in A_k, k = 1, 2, ..., n\}$ .

### **Measurable functions**

**DEFINITION 1.4 (measurable function – see Figure 1.1)** Given two measurable spaces  $(X, \mathscr{X})$  and  $(X', \mathscr{X}')$ , a function  $g : X \to X'$  is measurable if, for all  $A' \in \mathscr{X}'$ , the inverse image of A' through g, that is  $g^{-1}(A') := \{x \in X : g(x) \in A'\}$ , belongs to  $\mathscr{X}$ .

Sometimes, we emphasize one or both  $\sigma$ -algebras by writing  $\mathscr{X}$ -measurable or  $\mathscr{X}/\mathscr{X}'$ -measurable.

The importance of measurability becomes apparent when speaking of measures and measure spaces, as we shall do in the next section.

The next two theorems study the measurability of functions constructed from other measurable functions.

**THEOREM 1.5 (composition of measurable functions)** Given three measurable spaces  $(X, \mathscr{X})$ ,  $(X', \mathscr{X}')$ , and  $(X'', \mathscr{X}'')$ , and two measurable functions  $g: X \to X'$  and  $h: X' \to X''$ , the composition of g and h, i.e. the function  $h \cdot g: X \to X''$  defined through relation  $h \cdot g(x) := h(g(x))$ , is a  $\mathscr{X} / \mathscr{X}''$ -measurable function.

4



Figure 1.1: Measurable function.

PROOF. For any  $A'' \in \mathscr{X}''$ , we have

$$(h \cdot g)^{-1}(A'') = g^{-1}(h^{-1}(A'')).$$
(1.1)

Since *h* is  $\mathscr{X}'/\mathscr{X}''$ -measurable,  $h^{-1}(A'') \in \mathscr{X}'$ ; in turn, the  $\mathscr{X}/\mathscr{X}'$ -measurability of *g* implies that  $g^{-1}(h^{-1}(A'')) \in \mathscr{X}$ , which shows the  $\mathscr{X}/\mathscr{X}''$ -measurability of  $h \cdot g$ .  $\Box$ 

#### **THEOREM 1.6 (product and marginal measurable function)**

Consider the measurable space  $(X, \mathscr{X})$  and two other measurable spaces  $(X_1, \mathscr{X}_1)$ , and  $(X_2, \mathscr{X}_2)$  along with their product  $(X_1 \times X_2, \mathscr{X}_1 \otimes \mathscr{X}_2)$ .

**i)** Given two measurable functions  $g_1 : X \to X_1$  and  $g_2 : X \to X_2$ , the function  $g : X \to X_1 \times X_2$  defined according to the relation  $g(x) = (g_1(x), g_2(x)), x \in X$ , is a  $\mathscr{X} / \mathscr{X}_1 \otimes \mathscr{X}_2$ -measurable function from X to  $X_1 \times X_2$ ;

**ii**) conversely, if  $g: X \to X_1 \times X_2$  is  $\mathscr{X} / \mathscr{X}_1 \otimes \mathscr{X}_2$ -measurable, then  $g_1: X \to X_1$  such that  $g_1(x)$  is the first component of g(x) and the similarly defined  $g_2: X \to X_2$  are measurable functions from X to  $X_1$  and from X to  $X_2$ , respectively.

The theorem extends in an obvious way to the product of more measurable spaces.

### PROOF.

i) We need to show that  $g^{-1}(A) \in \mathscr{X}, \forall A \in \mathscr{X}_1 \otimes \mathscr{X}_2$ .

Consider the collection  $\mathscr{D}$  of all sets  $A \subseteq X_1 \times X_2$  such that  $g^{-1}(A) \in \mathscr{X}$ . We prove the following two facts:

(a)  $\mathscr{D}$  contains all sets of the form  $A = A_1 \times A_2, A_1 \in \mathscr{X}_1$  and  $A_2 \in \mathscr{X}_2$ ;

(b)  $\mathcal{D}$  is a  $\sigma$ -algebra.

Facts (a) and (b) imply the theorem thesis. Indeed, since  $\mathscr{X}_1 \otimes \mathscr{X}_2$  is the smallest  $\sigma$ -algebra that contains the sets of the form  $A_1 \times A_2$ ,  $A_1 \in \mathscr{X}_1$  and  $A_2 \in \mathscr{X}_2$ , it follows from (a) and (b) that  $\mathscr{X}_1 \otimes \mathscr{X}_2 \subseteq \mathscr{D}$ , so that any set in  $\mathscr{X}_1 \otimes \mathscr{X}_2$  has an inverse image through  $g^{-1}$  which is in  $\mathscr{X}$ .

(a) and (b) are proven as follows:

(a)  $g^{-1}(A_1 \times A_2) = g^{-1}((A_1 \times X_2) \cap (X_1 \times A_2)) = g^{-1}(A_1 \times X_2) \cap g^{-1}(X_1 \times A_2) = g_1^{-1}(A_1) \cap g_2^{-1}(A_2) \in \mathscr{X};$ (b) if  $A = \bigcup_{k=1}^{\infty} A_k$  with  $A_k \in \mathscr{D}$ , k = 1, 2, ..., then  $g^{-1}(A) = g^{-1}(\bigcup_{k=1}^{\infty} A_k) = \bigcup_{k=1}^{\infty} g^{-1}(A_k) \in \mathscr{X}$ , so that  $\mathscr{D}$  is closed under union. The fact that  $\mathscr{D}$  is closed under complementation and contains the entire set  $X_1 \times X_2$  is proven in a similar way.

ii) For any  $A_1 \in \mathscr{X}_1$ , we have that  $g_1^{-1}(A_1) = g^{-1}(A_1 \times X_2) \in \mathscr{X}$ , so that  $g_1$  is  $\mathscr{X}/\mathscr{X}_1$ -measurable. Similarly,  $g_2$  is  $\mathscr{X}/\mathscr{X}_2$ -measurable.

# **1.2** Measures and measure spaces

A measure is a function that associates to any set belonging to a  $\sigma$ -algebra a nonnegative number, the measure of the set. It must satisfy a certain set-theoretic property called  $\sigma$ -additivity.

**DEFINITION 1.7 (measure)** Let  $\mathscr{X}$  be a  $\sigma$ -algebra. A function  $m : \mathscr{X} \to [0,\infty]$  is called a measure if, for any countable collection of pairwise disjoint sets  $A_k \in \mathscr{X}$ , k = 1, 2, ..., the following property ( $\sigma$ -additivity) holds

$$m(\bigcup_{k=1}^{\infty} A_k) = \sum_{k=1}^{\infty} m(A_k).$$
 (1.2)

Definition 1.7 forces  $m(\emptyset) = 0$  (the empty set has measure zero). To see this, take  $A_1 = A$  and  $A_2 = A_3 = \cdots = \emptyset$  in (1.2). Moreover, by taking  $A_{k+1} = A_{k+2} = \ldots = \emptyset$  in (1.2) we see that the measure of the union of a finite number of sets equals the sum of their measures.

**DEFINITION 1.8 (measure space)** A triple  $(X, \mathcal{X}, m)$ , where X is any set,  $\mathcal{X}$  is a  $\sigma$ -algebra of subsets of X, and m is a measure on  $\mathcal{X}$  is called a measure space.  $\Box$ 

#### Image measures

Consider two measurable spaces  $(X, \mathscr{X})$  and  $(X', \mathscr{X}')$  and a measurable function  $g : X \to X'$ . If the first space is endowed with a measure m, i.e.  $(X, \mathscr{X}, m)$  is a measure space, then function g permits to define a measure m' on  $\mathscr{X}'$  according to the following definition:

$$m'(A') := m(g^{-1}(A')), A' \in \mathscr{X}'.$$
 (1.3)

(Here, one should note the importance of the fact that g is measurable: if g were not measurable, then  $g^{-1}(A')$  would need not be a set of  $\mathscr{X}$  so that  $m(g^{-1}(A'))$  could be undefined.) Measure m' is named the image measure of m through g.

#### Methods for introducing measures on measurable spaces

When defining a measure on a  $\sigma$ -algebra, it is sometimes convenient to first define the measure on a simpler system of sets and then to extend it to the  $\sigma$ -algebra.

A typical case is when the simpler system is an algebra  $\mathscr{A}$  (an algebra is a system of sets with the same properties as for a  $\sigma$ -algebra - see Definition 1.1 - where, however, property (b) is only required to hold for a finite number of sets) and the  $\sigma$ -algebra is  $\sigma(\mathscr{A})$ , the smallest  $\sigma$ -algebra containing  $\mathscr{A}$ . This situation is studied in the fundamental Theorem 1.9 below.

Before stating the theorem, we need some terminology. Given an algebra  $\mathscr{A}$ , a function  $m_0 : \mathscr{A} \to [0, \infty]$  satisfying (1.2) for all (possibly countably infinite) collection of pairwise disjoint sets  $A_k \in \mathscr{A}$ , k = 1, 2, ..., is called a premeasure (namely, a premeasure has identical properties as a measure but it is defined over an algebra instead of a  $\sigma$ -algebra). The reason for calling it a premeasure is that it extends naturally to a measure, as Theorem 1.9 states. A premeasure on  $\mathscr{A}$  is called  $\sigma$ -finite if there exits a sequence of sets  $A_k \in \mathscr{A}$ , k = 1, 2, ..., such that  $\bigcup_{k=1}^{\infty} A_k = X$  (i.e. the entire set) and  $m_0(A_k) < \infty, \forall k$ .

**THEOREM 1.9 (Caratheodory's)** Consider a  $\sigma$ -finite premeasure  $m_0$  on an algebra  $\mathscr{A}$ . Then, there is one and only one measure m that extends  $m_0$  to  $\sigma(\mathscr{A})$ , that is, a measure on  $\sigma(\mathscr{A})$  such that

$$m(A) = m_0(A), \quad for A \in \mathscr{A}.$$
 (1.4)

A proof can be found e.g. in the texts [3], [5] and [2, Theorem 3.1] for the specific case of probability measures.

As an application of Caratheodory's theorem, we next introduce the notion of product measure space.

#### **Product measure spaces**

Let us consider the product  $(X_1 \times X_2 \times \cdots \times X_n, \mathscr{X}_1 \otimes \mathscr{X}_2 \otimes \cdots \otimes \mathscr{X}_n)$  of *n* measurable spaces  $(X_k, \mathscr{X}_k), k = 1, 2, ..., n$ . We assume that each space  $(X_k, \mathscr{X}_k)$  is endowed with a  $\sigma$ -finite measure  $m_k$  and we want to introduce a product measure  $m_1 \times m_2 \times \cdots \times m_n$ on the product space.

For notational convenience, assume n = 2, the extension to n > 2 is straightforward. Consider the system of subsets of  $X_1 \times X_2$  of the form  $A_1 \times A_2$ , with  $A_1 \in \mathscr{X}_1$  and  $A_2 \in \mathscr{X}_2$ . Such a system is not an algebra since it is not closed under union. However, it is easily seen that the system of subsets consisting of finite unions of disjoint sets of the form  $A_1 \times A_2$  is indeed an algebra (though, not a  $\sigma$ -algebra.) For each set  $A = \bigcup_{k=1}^p (A_1^k \times A_2^k)$  of this algebra we define  $m_0(A) = \sum_{k=1}^p m_1(A_1^k)m_2(A_2^k)$ . It is a simple (but cumbersome) exercise to show that  $m_0$  is a  $\sigma$ -finite premeasure. Thus, by Caratheodory's extension Theorem 1.9,  $m_0$  can be extended in a unique way to a measure on the  $\sigma$ -algebra  $\mathscr{X}_1 \otimes \mathscr{X}_2$ . This measure is by definition the product measure and is denoted by  $m_1 \times m_2$ .

**DEFINITION 1.10 (product measure space)** Given *n* measure spaces  $(X_k, \mathscr{X}_k, m_k), k = 1, 2, ..., n$ , the measure space  $(X_1 \times X_2 \times \cdots \times X_n, \mathscr{X}_1 \otimes \mathscr{X}_2 \otimes \cdots \otimes \mathscr{X}_n, m_1 \times m_2 \times \cdots \times m_k)$  is called their product measure space.

## The measure space $(\mathbb{R}^n, \mathscr{B}(\mathbb{R}^n), \lambda^n)$

We start by considering the measure space  $(\mathbb{R}, \mathscr{B}(\mathbb{R}), \lambda)$ , i.e. we take n = 1, a space that plays a prominent role in measure theory.

Here,  $\mathbb{R}$  is the set of real numbers and  $\mathscr{B}(\mathbb{R})$  is the  $\sigma$ -algebra generated by all open intervals (a,b).  $\mathscr{B}(\mathbb{R})$  is named the Borel  $\sigma$ -algebra on the real line. Measure  $\lambda$  is the Lebesgue measure on the real line and is defined as follows. Consider the intervals of the form  $(a,b] := \{x \in \mathbb{R} : a < x \le b\}$ , or  $(-\infty,b] := \{x \in \mathbb{R} : x \le b\}$ , or  $(a,\infty) :=$  $\{x \in \mathbb{R} : a < x\}$ . Next, consider the system  $\mathscr{A}$  of subsets A of  $\mathbb{R}$  that are finite unions of disjoint intervals  $A_k$ , where each  $A_k$  has one of the indicated forms:  $A = \bigcup_{k=1}^p A_k$ .  $\mathscr{A}$  is an algebra. Define  $\lambda_0(A) := \sum_{k=1}^p (b_k - a_k)$ , where  $a_k$ ,  $b_k$  are the extremes of interval  $A_k$ . It can be seen that  $\lambda_0$  is a  $\sigma$ -finite premeasure, so that, by Caratheodory's Theorem 1.9, it can be extended in a unique way to a measure on  $\sigma(\mathscr{A})$ . Finally, it is an easy fact to prove that  $\sigma(\mathscr{A}) = \mathscr{B}(\mathbb{R})$ , so that the above construction has generated a measure on  $\mathscr{B}(\mathbb{R})$ , to which the name of Lebesgue measure  $\lambda$  is given.

**NOTE:** In many textbooks, the Lebesgue measure  $\lambda$  is defined over an extended  $\sigma$ -algebra obtained by augmenting  $\mathscr{B}(\mathbb{R})$  with all subsets of sets in  $\mathscr{B}(\mathbb{R})$  with zero  $\lambda$  measure. The  $\lambda$  measure of all these subsets is set to the zero value. This extended

 $\sigma$ -algebra is said to be "complete", an adjective that refers to the circumstance that any subset of a zero measure set is always measurable.

For  $n \ge 2$ ,  $(\mathbb{R}^n, \mathscr{B}(\mathbb{R}^n), \lambda^n)$  is simply defined as the *n*-fold product measure space of  $(\mathbb{R}, \mathscr{B}(\mathbb{R}), \lambda)$ .

Sometimes, it is of interest to consider the restriction of  $(\mathbb{R}, \mathscr{B}(\mathbb{R}), \lambda)$  to a set in  $\mathscr{B}(\mathbb{R})$ . An example is  $([0,1], \mathscr{B}[0,1], \lambda)$ . Here, by definition,  $\mathscr{B}[0,1]$  is the  $\sigma$ -algebra of all sets of the form  $A \cap [0,1]$ ,  $A \in \mathscr{B}(\mathbb{R})$  and  $\lambda$  is the restriction of the Lebesgue measure to these sets.

# **1.3 Integration**

Given a measure space  $(X, \mathcal{X}, m)$  and a  $\mathcal{X}/\mathcal{B}(\mathbb{R})$ -measurable function  $g : X \to \mathbb{R}$ , we want to define the integral of *g* with respect to the measure *m*, for which we use the symbol

$$\int_X g(x)dm(x).$$
(1.5)

This is done in 3 steps. For the first step we need the following definition.

**DEFINITION 1.11 (simple measurable function)** *Given a measurable space*  $(X, \mathscr{X})$ , *a measurable function*  $g : X \to \mathbb{R}$  *is said to be simple if it has the form* 

$$g = \sum_{k=1}^{N} \alpha_k \cdot 1(A_k), \quad A_k \in \mathscr{X}, \quad \alpha_k \in \mathbb{R}, \quad k = 1, 2, \dots, n,$$
(1.6)

where N is finite and  $1(A_k)$  is the indicator function of set  $A_k$ .

#### **STEP 1: Integral of non-negative simple measurable functions.**

Consider function g of the form in (1.6) and assume  $\alpha_k \ge 0, k = 1, 2, ..., n$ . Then, we define  $\int_X g(x) dm(x) = \sum_{k=1}^N \alpha_k m(A_k)$  (if  $\alpha_k = 0$  and  $m(A_k) = \infty$ , we let  $\alpha_k m(\alpha_k) = 0 \cdot \infty = 0$ ).

#### STEP 2: Integral of non-negative measurable functions.

Let now  $g: X \to \mathbb{R}$  be measurable and  $g \ge 0$ . Consider a sequence of simple measurable functions  $g_n$  such that  $g_n(x) \uparrow g(x), \forall x$  (that is, for all  $x, g_n(x)$  tends to g(x) and is increasing with n; in other words, it converges from below.) One such sequence is the following:

$$g_n := \sum_{k=1}^{n2^n} \frac{k-1}{2^n} \cdot 1(A_{n,k}) + n \cdot 1(A_n), \tag{1.7}$$

where  $A_{n,k} := \{x : (k-1)/2^n \le g(x) < k/2^n\}$  and  $A_n := \{x : g(x) \ge n\}$ . Clearly,  $\int_X g_n(x) dm(x)$  increases with *n*, so that it converges (either to a finite value or to  $+\infty$ .) We let  $\int_X g(x) dm(x) := \lim_{n\to\infty} \int_X g_n(x) dm(x)$ . One can prove that this definition is consistent, that is the limit is independent of the choice of the approximating sequence  $g_n$ .

## **STEP 3: Integral of measurable functions.**

Let  $g^+ := \max\{g, 0\}$  and  $g^- := -\min\{g, 0\}$  and note that  $g = g^+ - g^-$ . By definition, the integral of g is given by the formula

$$\int_{X} g(x) dm(x) = \int_{X} g^{+}(x) dm(x) - \int_{X} g^{-}(x) dm(x), \qquad (1.8)$$

provided that not both integrals in the right-hand side are  $+\infty$  (in which case we say that the integral is not defined).

### Notations

– When this generates no confusion, we drop the domain of integration and/or the arguments in the integral. So, for example, we write  $\int g dm$  for  $\int_X g(x) dm(x)$ .

– When the integral is performed over  $\mathbb{R}$  with respect to the Lebesgue measure  $\lambda$ , it is customary to write  $\int_{\mathbb{R}} g(x) dx$  for  $\int_{\mathbb{R}} g(x) d\lambda(x)$ .

- Take a set  $A \in \mathscr{X}$ . Then,  $g \cdot 1(A)$  (where 1(A) is the indicator function of set A) is measurable, provided that g is. We write  $\int_A g dm$  for  $\int_X g \cdot 1(A) dm$ . When A is the interval [a,b] and integration is with respect to the Lebesgue measure, we also write  $\int_a^b g(x) dx$ .

## Functions with value in $\overline{\mathbb{R}} = [-\infty, \infty]$

It is often of interest to integrate functions taking value on the extended real line  $[-\infty,\infty]$ . let  $\mathscr{B}(\bar{\mathbb{R}})$  be the  $\sigma$ -algebra generated by all intervals of the type (a,b),  $[-\infty,b)$ , and  $(a,\infty]$  and consider a  $\mathscr{X}/\mathscr{B}(\bar{\mathbb{R}})$ -measurable function g. The definition of integral extends to this case with no modifications with the agreement that  $0 \cdot \infty = 0$ .

#### **Properties of the integral**

The following properties of the integral are easy to prove (all integrals appearing in the formulas are assumed to exist by hypothesis).

- 1.  $\int 1(A)dm = m(A);$
- 2. if  $g_1 \leq g_2$ , then  $\int g_1 dm \leq \int g_2 dm$ ;

- 3.  $\int (\alpha g_1 + \beta g_2) dm = \alpha \int g_1 dm + \beta \int g_2 dm$ , provided that the right-hand side does not result in the indeterminate form  $\infty \infty$ ;
- 4. if g = 0 *m*-almost surely (i.e.  $m(\{g \neq 0\}) = 0$ ), then  $\int g dm = 0$ ;
- 5. if  $\int g \cdot 1(A) dm = 0, \forall A \in \mathscr{X}$ , then g = 0 *m*-almost surely.

### Change of space of integration

The following theorem permits to change space of integration in integrals.

#### **THEOREM 1.12** (change of space of integration in integrals)

Consider the three measurable spaces  $(X, \mathscr{X})$ ,  $(X', \mathscr{X}')$ , and  $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ , and two measurable functions  $g: X \to X'$  and  $h: X' \to \mathbb{R}$ . Further, assume that  $(X, \mathscr{X})$  is endowed with a measure m. Then,

$$\int_{X} h(g(x)) dm(x) = \int_{X'} h(x') dm'(x'), \qquad (1.9)$$

where m' is the image measure of m through g, in the sense that if one integral exists, then the other one also exists and the two are equal.

**PROOF.** For a simple non-negative function  $h = \sum_{k=1}^{N} \alpha_k \cdot 1(A_k)$ , we have

$$\int_{X} h(g(x)) dm(x) = \sum_{k=1}^{N} \alpha_{k} m(g^{-1}(A_{k}))$$
(1.10)

$$= \sum_{k=1}^{N} \alpha_k m'(A_k)$$
 (1.11)

$$= \int_{X'} h(x') dm'(x'). \qquad (1.12)$$

For a generic non-negative h, take a sequence of non-negative simple measurable functions  $h_n$  converging to h from below. Then,  $h_n \cdot g$  is a sequence of non-negative simple measurable functions defined on X converging from below to  $h \cdot g$  and, from what we have proven above, we obtain

$$\int_{X} h_n(g(x)) dm(x) = \int_{X'} h_n(x') dm'(x'), \quad \forall n.$$
(1.13)

Since the limit of the left-hand side as  $n \to \infty$  is by definition  $\int_X h(g(x)) dm(x)$  and that of the right-hand side is  $\int_{X'} h_n(x') dm'(x')$ , (1.9) is proven for non-negative *h*'s.

Finally, for a generic *h*, the result follows from the usual decomposition  $h = h^+ - h^-$ .

#### Integration with respect to a product measure

The following theorem permits to reduce the integral over  $X_1 \times X_2$  to an iterated integral (the proof can be found e.g. in [?], [5], [7], [2]).

**THEOREM 1.13 (Fubini's)** Let  $(X_1, \mathscr{X}_1, m_1)$  and  $(X_2, \mathscr{X}_2, m_2)$  be measure spaces with  $\sigma$ -finite measures  $m_1$  and  $m_2$  (a measure on  $(X, \mathscr{X})$  is  $\sigma$ -finite if X is the countable union of sets  $A_k \in \mathscr{X}$  with  $m(A_k) < \infty$ .) Further, let  $g: X_1 \times X_2 \to \mathbb{R}$  be a  $\mathscr{X}_1 \otimes \mathscr{X}_2$ -measurable function and assume that

$$\int_{X_1 \times X_2} |g| d(m_1 \times m_2) < \infty. \tag{1.14}$$

Then,

*i)*  $g(x_1, x_2)$  *is a*  $\mathscr{X}_1$ *-measurable function of*  $x_1$  *for each fixed*  $x_2$  *and a*  $\mathscr{X}_2$ *-measurable function of*  $x_2$  *for each fixed*  $x_1$ *;* 

ii) the integral  $\int_{X_1} g(x_1, x_2) dm_1$  is defined  $m_2$ -almost surely (i.e. the set where it is not defined is in  $\mathscr{X}_2$  and has  $m_2$  measure zero.) Moreover, if we define  $\int_{X_1} g(x_1, x_2) dm_1$  to be zero - or, equivalently, any other real number - where the integral is undefined, then  $\int_{X_1} g(x_1, x_2) dm_1$  is a  $\mathscr{X}_2$ -measurable function. A similar statement holds for  $\int_{X_2} g(x_1, x_2) dm_2$ ;

iii)

$$\int_{X_1 \times X_2} g(x_1, x_2) d(m_1 \times m_2) = \int_{X_1} \left[ \int_{X_2} g(x_1, x_2) dm_2 \right] dm_1$$
(1.15)

$$= \int_{X_2} \left[ \int_{X_1} g(x_1, x_2) dm_1 \right] dm_2, \qquad (1.16)$$

in the sense that the integral in the left-hand-side and the external integrals in the right-hand-side exist and equality holds.

The integral on the left is often referred to as the "double integral", while those on the right are called the "iterated integrals".  $\Box$ 

If we assume that g is nonnegative, the same result as in Fubini's theorem holds without requiring that the integral of |g| be bounded:

**THEOREM 1.14 (Tonelli's)** Let  $(X_1, \mathscr{X}_1, m_1)$  and  $(X_2, \mathscr{X}_2, m_2)$  be measure spaces with  $\sigma$ -finite measures  $m_1$  and  $m_2$  (a measure on  $(X, \mathscr{X})$ ) is  $\sigma$ -finite if X is

the countable union of sets  $A_k \in \mathscr{X}$  with  $m(A_k) < \infty$ .) Further, let  $g: X_1 \times X_2 \to \mathbb{R}$  be a  $\mathscr{X}_1 \otimes \mathscr{X}_2$ -measurable function with  $g \ge 0$ . Then,

i)  $g(x_1, x_2)$  is a  $\mathscr{X}_1$ -measurable function of  $x_1$  for each fixed  $x_2$  and a  $\mathscr{X}_2$ -measurable function of  $x_2$  for each fixed  $x_1$ ;

ii)  $\int_{X_1} g(x_1, x_2) dm_1$  is a  $\mathscr{X}_2$ -measurable function and  $\int_{X_2} g(x_1, x_2) dm_2$  is a  $\mathscr{X}_1$ -measurable function;

iii)

$$\int_{X_1 \times X_2} g(x_1, x_2) d(m_1 \times m_2) = \int_{X_1} \left[ \int_{X_2} g(x_1, x_2) dm_2 \right] dm_1$$
(1.17)

$$= \int_{X_2} \left[ \int_{X_1} g(x_1, x_2) dm_1 \right] dm_2.$$
 (1.18)

# Chapter 2

# **RANDOM VARIABLES**

In this Appendix, we make continuous reference to notions like measure, measurable space, and integration that are discussed in Appendix 1.

# 2.1 Random variables

**DEFINITION 2.1 (probability and probability space)** Given a measurable space  $(\Omega, \mathscr{F})$ , a probability  $\mathbb{P}$  is a measure on the  $\sigma$ -algebra  $\mathscr{F}$  such that  $\mathbb{P}(\Omega) = 1$ . The measure space  $(\Omega, \mathscr{F}, \mathbb{P})$  is called a probability space.

**DEFINITION 2.2 (random variable)** Given a probability space  $(\Omega, \mathscr{F}, \mathbb{P})$ , a  $\mathscr{F}/\mathscr{B}(\mathbb{R}^n)$ -measurable function  $v : \Omega \to \mathbb{R}^n$  is called a n-dimensional random variable. When n = 1, we simply speak of a random variable.  $\Box$ 

From Theorem 1.6, it is clear that n random variables form a n-dimensional random variable and viceversa.

## Interpretation and use of random variables

Random variables are used to describe phenomena in which a quantity of interest takes a value that remains unspecified at the moment the model is used, and yet one wants to incorporate in the model some beliefs on the chance with which the quantity will take value, beyond a set-theoretic description of its range of variability. According to Definition 2.2, corresponding to different points in  $\Omega$ , *v* assumes different real values and one can ask the question: what is the probability that *v* takes value in a given interval [a,b]? This probability is computable thanks to the assumption that v is  $\mathscr{F}/\mathscr{B}(\mathbb{R})$ measurable: the set of  $\omega \in \Omega$  such that  $v(\omega) \in [a,b]$  is an element of  $\mathscr{F}$  and so probability  $\mathbb{P}(\omega : v(\omega) \in [a,b])$  is provided by the model. Often, the value assumed by vis the only object of interest, while  $\Omega$  and  $\mathbb{P}$  are used as instruments to describe the probabilistic law with which different occurrences of the phenomenon take place.

### More on $\mathscr{B}(\mathbb{R}^n)$ and the measurability of random variables

By definition,  $\mathscr{B}(\mathbb{R})$  contains all open intervals (a,b). Since any open set on the real line is the countable union of open intervals, it is clear that  $\mathscr{B}(\mathbb{R})$  contains all open sets and it is in fact the  $\sigma$ -algebra generated by open sets (in measure theory, the term "Borel  $\sigma$ -algebra" is assigned to the  $\sigma$ -algebra generated by open sets in a given topological space. Thus,  $\mathscr{B}(\mathbb{R})$  is the Borel  $\sigma$ -algebra on  $\mathbb{R}$ ). Similarly,  $\mathscr{B}(\mathbb{R}^n)$  is the  $\sigma$ -algebra generated by the open sets in  $\mathbb{R}^n$ .

Suppose we want to prove that a given function  $v : \Omega \to \mathbb{R}^n$  is a random variable, i.e., it is  $\mathscr{F}/\mathscr{B}(\mathbb{R}^n)$ -measurable. In principle, we have to show that the inverse image of any set  $A \in \mathscr{B}(\mathbb{R}^n)$  is in  $\mathscr{F}$ . However, an easier test can be formulated: verify that the inverse image of just any open set is in  $\mathscr{F}$ . To see that this is enough, note that the system  $\mathscr{D}$  of all sets A in  $\mathbb{R}^n$  such that  $v^{-1}(A) \in \mathscr{F}$  is a  $\sigma$ -algebra (indeed, if  $A = \bigcup_{k=1}^{\infty} A_k$  with  $A_k \in \mathscr{D}$ , then  $v^{-1}(A) = v^{-1}(\bigcup_{k=1}^{\infty} A_k) = \bigcup_{k=1}^{\infty} v^{-1}(A_k) \in \mathscr{F}$ , so that  $\mathscr{D}$ is closed under union. The fact that  $\mathscr{D}$  is closed under complementation and contains the entire  $\mathbb{R}^n$  can be proven in a similar way). Now, since  $\mathscr{B}(\mathbb{R}^n)$  is the smallest  $\sigma$ algebra that contains the open sets, it is clear that  $\mathscr{B}(\mathbb{R}^n) \subseteq \mathscr{D}$ , showing that the inverse image of any set in  $\mathscr{B}(\mathbb{R}^n)$  is in  $\mathscr{F}$ , that is, the measurability of v.

Following the same reasoning, it is possible to conclude that *v* is measurable provided that the inverse image of any set in a system generating  $\mathscr{B}(\mathbb{R}^n)$  is in  $\mathscr{F}$ . For example, for n = 1 this leads to the following test:

**TEST OF MEASURABILITY 2.3**  $v : \Omega \to \mathbb{R}$  is measurable provided that  $v^{-1}(a,b) \in \mathscr{F}$  for any  $a, b \in \mathbb{R}$ .

Suppose now that  $f : \mathbb{R}^2 \to \mathbb{R}$  is continuous and  $v_1, v_2$  are two random variables. Then,  $f(v_1, v_2)$  is a random variable. So, e.g.  $v_1 + v_2, v_1 \cdot v_2, sin(v_1 \cdot v_2)$  are random variables. To see this, recall the definition of a continuous function: f is continuous if the inverse image of any open set is an open set. Thus, if f is continuous, it is  $\mathscr{B}(\mathbb{R}^2)/\mathscr{B}(\mathbb{R})$ -measurable. Appealing to Theorem 1.5 on the measurability of composition functions, we then conclude that  $f(v_1, v_2)$  is measurable. This fact extends in a natural way to functions  $f : \mathbb{R}^n \to \mathbb{R}$ .

## Probability distribution function and probability density function

To make notations easier, we consider first the case of 1-dimensional random variables. The extension to the n-dimensional case is discussed later in this section.

**DEFINITION 2.4 (probability distribution function)** *The probability distribution function (or, more simply, the distribution function) of a* 1-*dimensional random variable v is the function*  $F : \mathbb{R} \to [0, 1]$  *defined as*  $F(x) = \mathbb{P}(v^{-1}(-\infty, x])$ .

The following properties of F are a direct consequence of the properties of  $\mathbb{P}$  (the reader may want to try to detail a proof):

(a) F(x) is nondecreasing; (b) F(x) is continuous on the right; (c)  $\lim_{x\to-\infty} F(x) = 0$  and  $\lim_{x\to\infty} F(x) = 1$ .

If we let  $\mathbb{P}'$  be the image probability of  $\mathbb{P}$  on  $\mathscr{B}(\mathbb{R})$  induced by  $v(\mathbb{P}')$  is called the probability distribution or, more simply, the distribution), it is clear that  $F(x) = \mathbb{P}'(-\infty, x]$ , so that the probability distribution can be calculated from the image probability. It is an important fact that the converse is also true: the image probability  $\mathbb{P}'$  can be completely reconstructed from F. To see this, note that the system of subsets consisting of finite unions of disjoint sets of the form (a,b], or  $(-\infty,b]$ , or  $(a,\infty)$  is an algebra (let us call it  $\mathscr{A}$ ) and an element  $A := \bigcup_{k=1}^{n} A_k$  of this algebra has a probability that can be computed from F by the formula  $\mathbb{P}'(A) = \sum_{k=1}^{n} [F(b_k) - F(a_k)]$  (here,  $a_k$ ,  $b_k$  are the extremes of interval  $A_k$  and  $F(\infty)$  is short for  $\lim_{x\to\infty} F(x) = 1$  and similarly for  $F(-\infty)$ .) Then, by virtue of Caratheodory's Theorem 1.9, this probability can be extended in a unique way to  $\sigma(\mathscr{A}) = \mathscr{B}(\mathbb{R})$ , so reconstructing  $\mathbb{P}'$ .

Next, we define the probability density function.

**DEFINITION 2.5 (probability density function)** Suppose there exists a measurable function  $p : \mathbb{R} \to \mathbb{R}$  such that  $F(x) = \int_{-\infty}^{x} p(t)dt$ , where F is the probability distribution function of a random variable v. Then, p is called the probability density function (or, more simply, the density function) of the random variable v.  $\Box$ 

Given a random variable  $v : \Omega \to \mathbb{R}$  with distribution F and image probability P' and a measurable function  $g : \mathbb{R} \to \mathbb{R}$ , sometimes the notation  $\int_{\mathbb{R}} g(x)dF(x)$  is used in place of  $\int_{\mathbb{R}} g(x)d\mathbb{P}'(x)$ . It is easy to see that, if v admits density p, then  $\int_{\mathbb{R}} g(x)dF(x)$  is also equal to  $\int_{\mathbb{R}} g(x)p(x)dx$  in the sense that if one integral exists, also the other one exists and the two are equal (prove this by first considering simple non-negative g's, then non-negative g's and finally arbitrary measurable g's.)

From the definition of probability density function, it follows that, if p exists, then F is the integral of p and is therefore an absolute continuous function. But then F is  $\lambda$ -almost surely differentiable and p is  $\lambda$ -almost surely the derivative of F (this is

the "fundamental theorem of calculus", see e.g. Theorem 7.20 in [6]). Moreover, p is unique up to a set of zero Lebesgue measure, that is, if  $p_1$  and  $p_2$  are two densities associated with the same distribution F, then  $\lambda(x : p_1(x) \neq p_2(x)) = 0$ . We refer to different densities as "versions" of the density and a phrase like "consider the density of v" means: "consider a version of the density of v".

It is also worth mentioning that the density needs not exist. An example is given by a random variable taking always value 0. In this case

$$F(x) = \begin{cases} 0, & x < 0\\ 1, & x \ge 0, \end{cases}$$
(2.1)

and *p* clearly does not exist.

The above is an example of discrete distribution function. In general, we can identify three classes of distributions (discrete, absolutely continuous, and singular) and it turns out that any distribution is the convex combination of elements of these classes. This classification is made explicit in the following.

## **1. DISCRETE DISTRIBUTIONS**

A distribution function F is discrete if it is piecewise constant, that is it is constant except for certain points where it is discontinuous. The density p is not defined for discrete distributions.

By the following simple argument, it is possible to see that the number of points of discontinuity is countable, that is they can be enumerated as  $x_1, x_2, ...$ : there can be at most one point of discontinuity with jump bigger than 1/2 (with two jumps bigger than 1/2, *F* would reach a value bigger than 1); similarly, there can be at most three points of discontinuity with jump whose size belongs to (1/4, 1/2] and seven points with jump whose size belongs to (1/8, 1/4] and so on. Summing up, the points of discontinuity can be at most countable.

#### 2. ABSOLUTELY CONTINUOUS DISTRIBUTIONS

A distribution function F is absolutely continuous if it has probability density function.

Importantly, discrete and absolutely continuous distribution functions, or a combination of them, do not cover the set of all possibilities, that is, not all *F* can be written as  $F = \alpha F_d + (1 - \alpha)F_{ac}$ , for some  $\alpha \in [0, 1]$ , with  $F_d$  discrete and  $F_{ac}$  absolutely continuous. A universal decomposition is obtained by the introduction of a a third type of distributions, called singular:

$$F = \alpha F_d + \beta F_{ac} + (1 - \alpha - \beta) F_s, \qquad (2.2)$$

where  $F_s$  is a singular distribution function.

## **3. SINGULAR DISTRIBUTIONS**

A distribution function F is singular if it is continuous (and therefore any single point x has probability zero) but there exist a set of zero Lebesgue measure whose probability is 1.

Though singular distributions have a somehow peculiar behavior, examples of singular distributions are not difficult to construct, see e.g. [7]. Decomposition (2.2) is proven in many textbooks, among which [7].  $\Box$ 

#### More on probability distribution functions

Consider a probability measure  $\mathbb{P}$  on  $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ . Following the discussion after Definition 2.4 (in the discussion after Definition 2.4 we had the symbol  $\mathbb{P}'$  in place of  $\mathbb{P}$ ), we see that function  $F(x), x \in \mathbb{R}$ , defined as  $F(x) = \mathbb{P}(-\infty, x]$  satisfies properties (a), (b), (c), which we report here for the reader's convenience:

(a) F(x) is nondecreasing;
(b) F(x) is continuous on the right;
(c) lim<sub>x→-∞</sub> F(x) = 0 and lim<sub>x→∞</sub> F(x) = 1.

Moreover, P can be reconstructed from F.

Let us now ask a converse question: is it true that to an *F* satisfying (a), (b), (c), there always corresponds a (unique) probability  $\mathbb{P}$  such that  $\mathbb{P}(-\infty, x] = F(x)$ ? The answer is indeed positive.

**THEOREM 2.6** If F(x),  $x \in \mathbb{R}$ , satisfies (a), (b), (c), then there exists a unique probability measure  $\mathbb{P}$  on  $(\mathbb{R}, \mathscr{B}(\mathbb{R})$  such that  $P(-\infty, x] = F(x)$ ,  $\forall x \in \mathbb{R}$ .  $\Box$ 

Thus, there is a one-to-one correspondence between functions F(x) satisfying (a), (b), (c) and probability measures on  $(\mathbb{R}, \mathscr{B}(\mathbb{R}))$ .

Moreover, given a function F(x) satisfying (a), (b), (c) and the corresponding probability measure  $\mathbb{P}$ , the identity function v(x) = x is a random variable between  $(\mathbb{R}, \mathscr{B}(\mathbb{R}), \mathbb{P})$  and  $\mathbb{R}$  with distribution F(x). We thus see that any function F(x) satisfying (a), (b), (c) is a distribution function and conditions (a), (b), (c) characterize the set of all distribution functions.

**PROOF.** Consider the system of subsets of  $\mathbb{R}$  consisting of finite unions of disjoint sets of the form (a,b], or  $(-\infty,b]$ , or  $(a,\infty)$  and note that it is an algebra (let us call it  $\mathscr{A}$ ). To ease the notation, we write  $(a,\infty]$  for  $(a,\infty)$ , so that all intervals are written as (a,b] where *a* can possibly be  $-\infty$  and *b* can possibly be  $+\infty$ . For any set in  $\mathscr{A}$ , define

$$\mathbb{P}_0\left(\cup_{k=1}^p (a_k, b_k]\right) = \sum_{k=1}^p [F(b_k) - F(a_k)]$$
(2.3)

 $(F(\infty) \text{ is short for } \lim_{x\to\infty} F(x) = 1 \text{ and similarly for } F(-\infty))$ . If we prove that  $\mathbb{P}_0$  is countably additive, then by Caratheodory's Theorem 1.9 it can be extended in a unique way to a measure  $\mathbb{P}$  defined over  $\sigma(\mathscr{A}) = \mathscr{B}(\mathbb{R})$ . Moreover, for such a  $\mathbb{P}$  the following holds:

- $\mathbb{P}(-\infty, x] = F(x) 0 = F(x);$
- $\mathbb{P}(-\infty,\infty) = F(\infty) F(-\infty) = 1$ , showing that  $\mathbb{P}$  is a probability;
- no other probability P<sub>1</sub> ≠ P exists which satisfies P<sub>1</sub>(-∞,x] = F(x), x ∈ R. Indeed, the restriction of P<sub>1</sub> to A would satisfy (2.3) (with P<sub>1</sub> replacing P<sub>0</sub>); since P also satisfies (2.3) (with P replacing P<sub>0</sub>), by the uniqueness of the Caratheodory's extension we would have P<sub>1</sub> = P.

Thus, what is left to prove is that  $\mathbb{P}_0$  is countably additive on  $\mathscr{A}$ .

Let

$$A = \bigcup_{k=1}^{p} (a_k, b_k], \quad A_j = \bigcup_{k=1}^{p_j} (a_k^j, b_k^j], j = 1, 2, \dots,$$
(2.4)

where the  $A_j$ 's are disjoint and  $\bigcup_{j=1}^{\infty} A_j = A$ . Proving the countable additivity of  $\mathbb{P}_0$  amounts to show that

$$\mathbb{P}_0(A) = \sum_{j=1}^{\infty} \mathbb{P}_0(A_j).$$
(2.5)

Rewriting (2.5) as follows

$$0 = \mathbb{P}_0(A) - \sum_{j=1}^{\infty} \mathbb{P}_0(A_j) = \mathbb{P}_0(A) - \lim_{m \to \infty} \sum_{j=1}^{m} \mathbb{P}_0(A_j) = \lim_{m \to \infty} \mathbb{P}_0(A - \bigcup_{j=1}^{m} A_j) \quad (2.6)$$

and observing that  $A - \bigcup_{j=1}^{m} A_j =: B_m \downarrow \emptyset$  (" $\downarrow$ " means that  $B_m \supseteq B_{m+1}, m = 1, 2, ...,$ and  $\bigcap_{m=1}^{\infty} B_m = \emptyset$ , the empty set), we see that (2.5) can be rewritten as

$$\lim_{m\to\infty}\mathbb{P}_0(B_m)=0, \quad \text{for any sequence } B_m \in \mathscr{A}, m=1,2,\ldots, \text{ such that } B_m \downarrow \emptyset. (2.7)$$

The proof is now completed by showing the validity of (2.7).

Let us suppose first that  $B_m \in [-M, M]$ , m = 1, 2, ..., for some  $M < \infty$ . Since F(x) is continuous on the right, the left extremes of the intervals forming  $B_m$  can be slightly moved to the right without a significant change of the F value. Therefore, we can

find sets  $C_m \in \mathscr{A}$  such that  $closure(C_m) \subseteq B_m$  and  $\mathbb{P}_0(B_m - C_m) \leq \frac{1}{2^m}\varepsilon$ , where  $\varepsilon > 0$  is a preassigned small number. One fundamental property of the sets  $closure(C_m)$ , m = 1, 2, ..., is that the intersection of a finite number of them is already empty:

$$\bigcap_{m=1}^{r} closure(C_m) = \emptyset, \quad r < \infty.$$
(2.8)

The reason is that sets  $[-M,M] - closure(C_m)$ , m = 1, 2, ..., form an open covering of [-M,M]. But, [-M,M] is compact (by Heine-Borel theorem) so that a finite subcovering exists:

$$\cup_{m=1}^{r}([-M,M] - closure(C_m)) = [-M,M],$$
(2.9)

and this implies (2.8).

Thus,

$$\mathbb{P}_{0}(B_{r}) = \mathbb{P}_{0}(B_{r} - \bigcap_{m=1}^{r} C_{m}) \quad (use \ (2.8))$$
(2.10)

$$\leq \mathbb{P}_0(\cup_{m=1}^r (B_m - C_m)) \tag{2.11}$$

$$\leq \sum_{m=1}^{\prime} \mathbb{P}_0(B_m - C_m) \tag{2.12}$$

$$\leq \sum_{m=1}^{r} \frac{1}{2^m} \varepsilon \tag{2.13}$$

$$= \varepsilon,$$
 (2.14)

from which (2.7) follows.

Suppose now that sets  $B_m$  are not confined to an interval [-M,M]. Then, conclusion (2.7) can still be drawn by taking an interval [-M,M] such that F(-M) and 1 - F(M) are smaller than  $\varepsilon$  and then following the same line of reasoning as before (details are left to the reader).

#### Multidimensional random variables

The notions of distribution function and density function can be extended with just some notational complications to multidimensional random variables, while maintaining all the fundamental properties valid for the 1-dimensional case. Here, it suffices to say that the probability distribution function is defined as

$$F(x_1, x_2, \dots, x_n) = \mathbb{P}(v_1^{-1}(-\infty, x_1] \cap v_2^{-1}(-\infty, x_2] \cap \dots \cap v_n^{-1}(-\infty, x_n])$$
(2.15)  
=  $\mathbb{P}'((-\infty, x_1] \times (-\infty, x_2] \times \dots \times (-\infty, x_n]),$  (2.16)

where  $v_k$  is the k-th component of the *n*-dimensional random variable v and  $\mathbb{P}'$  is the image probability. Again, it is possible to see that *F* uniquely defines  $\mathbb{P}'$ .

In the multidimensional case, the probability distribution has the following properties that extend those valid for the 1-dimensional case:

(a) Given two points  $x^{(1)} = (x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})$  and  $x^{(2)} = (x_1^{(2)}, x_2^{(2)}, \dots, x_n^{(2)})$  in  $\mathbb{R}^n$ with  $x_1^{(2)} \ge x_1^{(1)}, x_2^{(2)} \ge x_2^{(1)}, \dots, x_n^{(2)} \ge x_n^{(1)}$ , and a function  $f : \mathbb{R}^n \to \mathbb{R}$ , introduce the notation  $\Delta_k^{x_k^{(1)}, x_k^{(2)}} f := f(x_1, \dots, x_{k-1}, x_k^{(2)}, x_{k+1}, \dots, x_n) - f(x_1, \dots, x_{k-1}, x_k^{(1)}, x_{k+1}, \dots, x_n)$ . Then, it is a matter of a cumbersome computation to show that  $\Delta_1^{x_1^{(1)}, x_1^{(2)}}(\Delta_2^{x_2^{(1)}, x_2^{(2)}}(\dots (\Delta_n^{x_n^{(1)}, x_n^{(2)}}F))) = P'((x_1^{(1)}, x_1^{(2)}] \times (x_2^{(1)}, x_2^{(2)}] \times \dots \times (x_n^{(1)}, x_n^{(2)}])$ . Since the right-hand side is clearly nonnegative, we then have

$$\Delta_1^{x_1^{(1)}, x_1^{(2)}} \cdots \Delta_n^{x_n^{(1)}, x_n^{(2)}} F \ge 0.$$
(2.17)

To help visualize the situation, for n = 2 we have  $\Delta_1^{x_1^{(1)}, x_1^{(2)}} (\Delta_2^{x_2^{(1)}, x_2^{(2)}} F) = F(x_1^{(2)}, x_1^{(2)}) - F(x_1^{(1)}, x_2^{(1)}) - F(x_1^{(2)}, x_2^{(1)}) + F(x_1^{(1)}, x_2^{(1)})$  and it represents the measure P' of the rectangle in Figure 2.1.



Figure 2.1:  $\Delta_1^{x_1^{(1)},x_1^{(2)}}(\Delta_2^{x_2^{(1)},x_2^{(2)}}F)$  is the measure *P*' of the rectangle.

For n = 1, (2.17) is equivalent to say that F(x) is nondecreasing; (b)  $F(x_1, x_2, ..., x_n)$  is continuous on the right in the sense that if  $x_k \downarrow \bar{x}_k, k = 1, 2, ..., n$ , then  $F(x_1, x_2, ..., x_n) \rightarrow F(\bar{x}_1, \bar{x}_2, ..., \bar{x}_n)$ ; (c) if  $x_k \rightarrow \bar{x}_k, k = 1, 2, ..., n$ , and at least one of the  $\bar{x}_k$ 's is  $-\infty$ , then  $F(x_1, x_2, ..., x_n) \rightarrow F(\bar{x}_1, \bar{x}_2, ..., \bar{x}_n)$ 

0. Moreover,  $\lim_{x_1\to\infty,\dots,x_n\to\infty} F(x_1,x_2,\dots,x_n) = 1$ .

A *n*-dimensional probability density function is a nonnegative measurable function  $\mathbb{P}$  such that  $F(x_1, x_2, ..., x_n) = \int_{-\infty}^{x_1} \left[ \int_{-\infty}^{x_2} \cdots \left[ \int_{-\infty}^{x_n} p(t_1, t_2, ..., t_n) dt_n \right] \cdots dt_2 \right] dt_1.$ 

#### **Expectation, variance and moments**

Given a random variable v, the integral  $\int_{\Omega} v d\mathbb{P}$  (assuming it is defined) is called the expectation (or mean) of v and is also written  $\mathbb{E}[v]$ . Other significant integral characteristics of v are its moment or order  $r : \mathbb{E}[v^r]$  (that is the expectation of the random variable obtained by the composition of v with the  $\mathbb{R} \to \mathbb{R}$  function of elevation to the *r*-th power), and its variance:  $\mathbb{E}[(v - \mathbb{E}[v])^2]$ . To indicate the variance of v, the symbol var(v) is also used.

When *v* is a matrix of random variables with entries  $v_{kj}$ , by the symbol  $\mathbb{E}[v]$  we mean the matrix with entries  $\mathbb{E}[v_{kj}]$ . If *v* is a *n*-dimensional random variable,  $\mathbb{E}[v]$  is the vector listing the expectation of its components. Similarly, var(v) is a matrix with entries  $\mathbb{E}[(v_j - \mathbb{E}[v_j])(v_k - \mathbb{E}[v_k])]$ , where  $v_j$ ,  $v_k$  are the components of *v*.

Note that the expectation and the other integral quantities can be computed in different ways. For example, letting  $F^{(v)}$  be the probability distribution of v and  $F^{(v^2)}$  that of  $v^2$ , we have

$$\mathbb{E}[v^2] = \int_{\Omega} v^2 dP = \int_{\mathbb{R}} x^2 dF^{(v)}(x) = \int_{\mathbb{R}} x dF^{(v^2)}(x), \qquad (2.18)$$

where the last two equalities are justified in the light of Theorem 1.12.

# 2.2 Independence and uncorrelation

We start by considering two 1-dimensional random variables  $v_1$  and  $v_2$ .

**DEFINITION 2.7 (independence)** We say that  $v_1$  and  $v_2$  are independent if

$$\mathbb{P}(v_1^{-1}(A_1) \cap v_2^{-1}(A_2)) = \mathbb{P}(v_1^{-1}(A_1)) \cdot \mathbb{P}(v_2^{-1}(A_2)), \quad \forall A_1, A_2 \in \mathscr{B}(\mathbb{R}).$$
(2.19)

Hence, if  $v_1$  and  $v_2$  are independent, then the probability that they simultaneously take value in given ranges  $A_1$  and  $A_2$  equals the product of the probabilities that the first one takes value in  $A_1$  and that the second one takes value in  $A_2$ .

If  $v_1$  and  $v_2$  are independent, so are  $f(v_1)$  and  $g(v_2)$ , with f and g are arbitrary measurable functions (show this).

Let  $\mathbb{P}'_1$  and  $\mathbb{P}'_2$  be the image probabilities on  $\mathbb{R}$  induced by two independent random variables  $v_1$  and  $v_2$ , respectively. Also, consider the 2-dimensional random variable defined through relation  $v = (v_1, v_2)$  and let  $\mathbb{P}'$  be the corresponding image probability on  $\mathbb{R}^2$ . It turns out that  $\mathbb{P}' = \mathbb{P}'_1 \times \mathbb{P}'_2$ . To prove this, it suffices to show that  $\mathbb{P}'$  and  $\mathbb{P}'_1 \times \mathbb{P}'_2$  agrees over the algebra of finite unions of disjoint sets of the form  $A_1 \times A_2$  with  $A_1, A_2 \in \mathscr{B}(\mathbb{R})$ . In fact, by Caratheodory's theorem 1.9 we then have that the extension to  $\mathscr{B}(\mathbb{R}^2)$  is unique and, therefore, coincident. Take  $A = \bigcup_{k=1}^p (A_1^k \times A_2^k)$ . We have:  $\mathbb{P}'(A) = \mathbb{P}'(\bigcup_{k=1}^p (A_1^k \times A_2^k)) = \sum_{k=1}^p \mathbb{P}'(A_1^k \times A_2^k) =$ [due to independence of  $v_1$  and  $v_2$ ] =  $\sum_{k=1}^p \mathbb{P}'_1(A_1^k) \cdot \mathbb{P}'_2(A_2^k) = (\mathbb{P}'_1 \times \mathbb{P}'_2)(A)$ , so that  $\mathbb{P}'$ and  $\mathbb{P}'_1 \times \mathbb{P}'_2$  indeed agree over the considered algebra.

Letting  $F_1$  and  $F_2$  be the probability distribution functions of  $v_1$  and  $v_2$  and F that of v, as a direct consequence of (2.19) we have that  $F(x_1, x_2) = F_1(x_1) \cdot F_2(x_2)$ . Moreover, if  $F_1$  and  $F_2$  admit density function, say  $p_1$  and  $p_2$ , we then have that v has density function too and it is given by  $p(x_1, x_2) = p_1(x_1) \cdot p_2(x_2)$ , as it is shown by direct inspection:

$$\int_{-\infty}^{x_1} \left[ \int_{-\infty}^{x_2} p_1(t_1) p_2(t_2) dt_2 \right] dt_1 = \int_{-\infty}^{x_1} p_1(t_1) \left[ \int_{-\infty}^{x_2} p_2(t_2) dt_2 \right] dt_1 \quad (2.20)$$

$$= \int_{-\infty}^{x_1} p_1(t_1) F_2(x_2) dt_1 \qquad (2.21)$$

$$= F_1(x_1) \cdot F_2(x_2) \tag{2.22}$$

$$= F(x_1, x_2). (2.23)$$

The converse also holds true: if  $F(x_1, x_2) = F_1(x_1) \cdot F_2(x_2)$ , or  $p(x_1, x_2) = p_1(x_1) \cdot p_2(x_2)$ , then  $v_1$  and  $v_2$  are independent (providing details is a useful exercise).

**DEFINITION 2.8 (uncorrelation)** We say that  $v_1$  and  $v_2$  are uncorrelated if  $\mathbb{E}[v_1v_2]$ ,  $\mathbb{E}[v_1]$  and  $\mathbb{E}[v_2]$  exist finite and

$$\mathbb{E}[v_1 v_2] = \mathbb{E}[v_1] \cdot \mathbb{E}[v_2]. \tag{2.24}$$

Uncorrelation is an integral notion. Not surprisingly, independence is a stronger notion than uncorrelation and the former implies the latter, while the opposite is in general false. More precisely, suppose that  $v_1$  and  $v_2$  are independent and that  $E[v_1v_2]$ ,  $E[v_1]$  and  $E[v_2]$  exist finite; then,  $v_1$  and  $v_2$  are uncorrelated, as the following calculation shows:

$$\mathbb{E}[v_1 v_2] = \int_{\Omega} v_1(\boldsymbol{\omega}) v_2(\boldsymbol{\omega}) d\mathbb{P}(\boldsymbol{\omega})$$
(2.25)

$$= \int_{\mathbb{R}^2} xy \, d\mathbb{P}'(x, y) \quad (use \ Theorem \ 1.12) \tag{2.26}$$

$$= \int_{\mathbb{R}^2} xy \, d(\mathbb{P}'_1 \times \mathbb{P}_2)(x, y) \tag{2.27}$$

$$= \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} xy \, d\mathbb{P}'_1(x) \right] d\mathbb{P}'_2(y) \quad (use \ Theorem \ 1.13) \qquad (2.28)$$

$$= \left[ \int_{\mathbb{R}} x \, d\mathbb{P}'_1(x) \right] \left[ \int_{\mathbb{R}} y \, d\mathbb{P}'_2(y) \right]$$
(2.29)

$$= \mathbb{E}[v_1]\mathbb{E}[v_2]. \tag{2.30}$$

An example of two random variables that are uncorrelated but not independent is shown in Figure 2.2.



Figure 2.2: Probability space:  $(\Omega, \mathscr{F}, \mathbb{P}) = ([0, 1], \mathscr{B}[0, 1], \lambda)$ . Solid line =  $v_1$ ; dashed line =  $v_2$ . Left:  $v_1$  and  $v_2$  are uncorrelated, but not independent; right:  $v_1$  and  $v_2$  are independent, and therefore also uncorrelated.

The notions of independence and uncorrelation carry over to the multidimensional case in a straightforward way. Given  $v_1 : \Omega \to \mathbb{R}^{n_1}$  and  $v_2 : \Omega \to \mathbb{R}^{n_2}$ , we say that  $v_1$  is independent of  $v_2$  if  $\mathbb{P}(v_1^{-1}(A_1) \cap v_2^{-1}(A_2)) = \mathbb{P}(v_1^{-1}(A_1)) \cdot \mathbb{P}(v_2^{-1}(A_2)), \forall A_1 \in \mathscr{B}(\mathbb{R}^{n_1}), A_2 \in \mathscr{B}(\mathbb{R}^{n_2})$ . Note that this definition only establish a cross-property of  $v_1$  and  $v_2$ ; different components of e.g.  $v_1$  can well be dependent one on the others. By identifying  $v_1$ and  $v_2$  with the vectors of their components, we say that  $v_1$  and  $v_2$  are uncorrelated if  $\mathbb{E}[v_1v_2^T] = \mathbb{E}[v_1]\mathbb{E}[v_2^T]$ .

# 2.3 Characteristic functions

The method of characteristic functions is one of the main tools in probability theory. Though a characteristic function carries exactly the same information content as the corresponding probability distribution, in many contexts it is more handy to use than the distribution itself. Here, we merely define characteristic functions and derive some basic properties of use in the book. The interested reader is referred to textbooks on probability for a broader treatment.

To define the characteristic function, we need to use complex-valued random variables. A complex-valued random variable v is, by definition, given by  $v = v_{\mathbb{R}} + iv_{\mathbb{I}}$ , where  $v_{\mathbb{R}}$  and  $v_{\mathbb{I}}$  are (real-valued) random variables. We also let  $\mathbb{E}[v] = \mathbb{E}[v_{\mathbb{R}}] + i\mathbb{E}[v_{\mathbb{I}}]$ .

**DEFINITION 2.9 (characteristic function)** *The characteristic function of a random variable v is defined as*  $\varphi(t) := \mathbb{E}[e^{itv}], t \in \mathbb{R}$ .

For a given t,  $\mathbb{E}[e^{itv}]$  is a complex number; as t varies over  $\mathbb{R}$ ,  $\varphi(t) = \mathbb{E}[e^{itv}]$  is a complex-valued function. It is clear that  $\varphi(t)$  can also be expressed as  $\varphi(t) = \int_{\mathbb{R}} e^{itx} dF(x)$ , where F is the distribution function of v. Thus,  $\varphi$  is determined by F. It is a crucial fact that the converse is also true: F can be completely reconstructed from  $\varphi$ , as the next theorem states.

**THEOREM 2.10** Let *F* and *G* be probability distribution functions on  $\mathbb{R}$  with the same characteristic function, viz.,

$$\int_{\mathbb{R}} e^{itx} dF(x) = \int_{\mathbb{R}} e^{itx} dG(x), \quad \forall t \in \mathbb{R}.$$
(2.31)

Then, F(x) = G(x),  $\forall x \in \mathbb{R}$ .

**PROOF.** Consider the function  $f^{\varepsilon}$  in Figure 2.3, where  $\beta > \alpha$  are arbitrary and  $\varepsilon > 0$  is smaller than  $\beta - \alpha$ .



Figure 2.3:

We first prove that

$$\int_{\mathbb{R}} f^{\varepsilon} dF = \int_{\mathbb{R}} f^{\varepsilon} dG, \qquad (2.32)$$

and then show that the thesis of the theorem follows from this equality.

Consider a sequence of real numbers  $\rho_n \downarrow 0$  and pick a *n* large enough so that  $[\alpha, \beta + \varepsilon] \subseteq [-n, n]$ . In [-n, n],  $f^{\varepsilon}$  can be uniformly approximated by a finite trigonometric sum (Weierstrass theorem). Precisely, there exists a function

$$f_n^{\varepsilon}(x) := \sum_{k=-N(n)}^{N(n)} a_k e^{i\pi\frac{k}{n}x},$$
(2.33)

where  $a_k$  are complex coefficients, such that

$$\sup_{-n \le x \le n} |f^{\varepsilon}(x) - f_n^{\varepsilon}(x)| \le \rho_n.$$
(2.34)

Observe also that function  $f_n^{\varepsilon}(x)$  is periodic so that

$$\sup_{x} |f_n^{\varepsilon}(x)| = \sup_{-n \le x \le n} |f_n^{\varepsilon}(x)| \le 1 + \rho_n,$$
(2.35)

and that, by (2.31),

$$\int_{\mathbb{R}} f_n^{\varepsilon} dF = \int_{\mathbb{R}} f_n^{\varepsilon} dG.$$
(2.36)

Thus,

$$\left| \int_{\mathbb{R}} f^{\varepsilon} dF - \int_{\mathbb{R}} f^{\varepsilon} dG \right|$$
(2.37)

$$= \left| \int_{[-n,n]} f^{\varepsilon} dF - \int_{[-n,n]} f^{\varepsilon} dG \right|$$
(2.38)

$$\leq \left| \int_{[-n,n]} f_n^{\varepsilon} dF - \int_{[-n,n]} f_n^{\varepsilon} dG \right| + 2\rho_n \quad (use \ (2.34))$$

$$(2.39)$$

$$\leq \left| \int_{\mathbb{R}} f_n^{\varepsilon} dF - \int_{\mathbb{R}} f_n^{\varepsilon} dG \right| + (1+\rho_n) \int_{[-n,n]^c} dF + (1+\rho_n) \int_{[-n,n]^c} dG + 2\rho_{\sigma}^2 .40)$$

$$(use \ (2.35); let [-n,n]^c be the complement \ of \ [-n,n] (2.41)$$

$$\leq (1+\rho_n) \int_{[-n,n]^c} dF + (1+\rho_n) \int_{[-n,n]^c} dG + 2\rho_n \quad (use \ (2.36)). \tag{2.42}$$

The right-hand side tends to zero as  $n \to \infty$ . Since the left-hand side does not depend on *n*, it must then be equal to zero and (2.32) is proven.

We turn now to prove that (2.32) implies that F = G. As  $\varepsilon \to 0$ ,  $\int_{\mathbb{R}} f^{\varepsilon} dF \to F(\beta) - F(\alpha)$  and  $\int_{\mathbb{R}} f^{\varepsilon} dG \to G(\beta) - G(\alpha)$ . Hence, from (2.32),  $F(\beta) - F(\alpha) = G(\beta) - G$ 

$$G(\alpha)$$
. Letting  $\alpha \to -\infty$ , we conclude that  $F(\beta) = G(\beta), \forall \beta \in \mathbb{R}$ .

Why are characteristic functions so widely used in probability? One important reason is that the characteristic function of the sum of independent random variables is simply given by the product of the characteristic functions of the variables:

- If  $v_1$  and  $v_2$  are independent, then  $\varphi_{v_1+v_2}(t) = \varphi_{v_1}(t) \cdot \varphi_{v_2}(t)$ .

To see that this is the case, write:  $\varphi_{v_1+v_2}(t) = \mathbb{E}[e^{it(v_1+v_2)}] = \mathbb{E}[e^{itv_1}e^{itv_2}] =$ [since  $v_1$  and  $v_2$  are independent, so are  $e^{itv_1}$  and  $e^{itv_2}$  and independence implies uncorrelation] =  $\mathbb{E}[e^{itv_1}]\mathbb{E}[e^{itv_2}] = \varphi_{v_1}(t) \cdot \varphi_{v_2}(t)$ . So, when dealing with independent variables, we can move from distributions (for which independence translates into the awkward condition that the distribution of  $v_1 + v_2$  is the convolution of the distributions of  $v_1$  and  $v_2$ ) to characteristic functions and use the handy product rule. In doing so, no information is lost, as the distribution can be reconstructed from the characteristic function, as stated in Theorem 2.10.

We know that  $\varphi_1(t) = \varphi_2(t)$  implies  $F_1(x) = F_2(x)$ . Now, we ask: suppose that  $\varphi_n(t) \rightarrow \varphi(t)$ ; is it true that  $F_n(x) \rightarrow F(x)$ ? A precise answer is given by the following theorem, which plays an important role in proving limit results in probability theory.

**THEOREM 2.11** Let  $F_n$  be a sequence of probability distribution functions on  $\mathbb{R}$  and let  $\varphi_n$  be the corresponding sequence of characteristic functions.

- (a) If  $F_n \to F$  weakly (see Section 3.5 for the notion of weak convergence) and  $\varphi$  is the characteristic function of F, then  $\varphi_n(t) \to \varphi(t)$ ,  $\forall t \in \mathbb{R}$ ;
- (b) if  $\varphi_n(t) \to \varphi(t)$ ,  $\forall t \in \mathbb{R}$ , and  $\varphi(t)$  is continuous at t = 0, then  $\varphi(t)$  is a characteristic function (i.e.,  $\varphi(t) = \int_{\mathbb{R}} e^{itx} dF(x)$  for some distribution function F) and  $F_n \to F$  weakly.

#### PROOF.

(a) Write  $\varphi_n(t) = \int_{\mathbb{R}} (\cos(tx) + i\sin(tx)) dF_n(x)$ . The weak convergence of  $F_n \to F$  means that  $\int_{\mathbb{R}} f(x) dF_n(x) \to \int_{\mathbb{R}} f(x) dF(x)$ , for any continuous and bounded function f(x). The thesis then follows by taking in turn  $f(x) = \cos(tx)$  and  $f(x) = \sin(tx)$ .

(b) The proof proceeds as follows. Thanks to the continuity of  $\varphi(t)$  at t = 0, we prove that  $F_n$  is tight (see Theorem 3.24 for the definition of tightness). Due to tightness, by Theorem 3.24,  $F_n$  admits a subsequence weakly convergent to some F and this F has  $\varphi(t)$  as characteristic function. Finally, by the convergence  $\varphi_n(t) \rightarrow \varphi(t)$  we establish that the whole sequence  $F_n$  converges to F.

=

To prove the tightness of  $F_n$ , pick any real M and let  $\beta := \inf_{|\alpha| \ge 1} \left(1 - \frac{\sin \alpha}{\alpha}\right)$  (note that  $\beta > 0$ ). We have

$$\beta \int_{|x| \ge M} dF_n(x) = \inf_{|\alpha| \ge 1} \left( 1 - \frac{\sin \alpha}{\alpha} \right) \int_{|x| \ge M} dF_n(x)$$
(2.43)

$$\leq \int_{|x|\geq M} \left(1 - \frac{\sin(x/M)}{x/M}\right) dF_n(x) \tag{2.44}$$

$$\leq \int_{\mathbb{R}} \left( 1 - \frac{\sin(x/M)}{x/M} \right) dF_n(x)$$
(2.45)

$$(for \ x = 0, \ let \ \frac{\sin(x/M)}{x/M} \ be \ 1)$$
 (2.46)

$$= \int_{\mathbb{R}} \left[ M \int_{0}^{1/M} (1 - \cos(tx)) dt \right] dF_n(x)$$
(2.47)

$$M \int_0^{1/M} \left[ \int_{\mathbb{R}} (1 - \cos(tx)) dF_n(x) \right] dt \qquad (2.48)$$

$$(use Fubini's Theorem 1.13)$$
(2.49)

$$= M \int_{0}^{1/M} (1 - Re(\varphi_n(t))) dt$$
 (2.50)

$$\stackrel{n \to \infty}{\longrightarrow} \quad M \int_0^{1/M} (1 - Re(\varphi(t))) dt$$
 (2.51)

(use a slight variation of the dominated (2.52)

$$convergence \ Theorem \ 3.9). \tag{2.53}$$

The right-hand side represents the mean value of  $1 - Re(\varphi(t))$  in a right neighborhood of the origin. Since  $\varphi(0) = \mathbb{E}[e^{i0x}] = 1$  and  $\varphi(t)$  is continuous at t = 0, we have

$$M \int_0^{1/M} (1 - Re(\varphi(t))) dt \to 0, \quad as \ M \to \infty.$$
(2.54)

We show that (2.54) implies the tightness of  $F_n$ . Indeed, given an arbitrarily small  $\varepsilon > 0$ , take  $M(\varepsilon)$  such that  $M(\varepsilon) \int_0^{1/M(\varepsilon)} (1 - Re(\varphi(t))) dt \le \frac{\varepsilon}{2}$ . Then, by (2.53),  $\beta \int_{|x| \ge M(\varepsilon)} dF_n(x) \le \varepsilon$  for any *n* large enough, say  $n \ge n(\varepsilon)$ . Since  $\beta \int_{|x| \ge M} dF_n(x)$  is decreasing with *M*, we then have  $\sup_{n \ge n(\varepsilon)} \beta \int_{|x| \ge M} dF_n(x) \le \varepsilon$ , as  $M \to \infty$ . On the other hand, over the finite set of integers *n* with  $n < n(\varepsilon)$ , we have  $\max_{n < n(\varepsilon)} \beta \int_{|x| \ge M} dF_n(x) \ge \varepsilon$ , as  $M \to \infty$ . Putting together these two facts yields:  $\sup_n \beta \int_{|x| \ge M} dF_n(x) \le \varepsilon$ , as  $M \to \infty$ . Owing to the arbitrariness of  $\varepsilon$ , the tightness of  $F_n$  follows.

Having proven the tightness of  $F_n$ , appeal now to Helly's Theorem 3.24 to conclude that there exists a subsequence  $F_{n_k}$  of  $F_n$  that converges weakly to some limit distribution function F. Since  $F_{n_k} \to F$  weakly, from part (a) of this theorem, we also have

that  $\varphi_{n_k}(t)$  tends for every *t* to the characteristic function of *F*. But, by assumption,  $\varphi_{n_k}(t) \rightarrow \varphi(t), \forall t \in \mathbb{R}$ , so that  $\varphi(t)$  must be the characteristic function of *F*.

We conclude the proof by showing that the whole sequence  $F_n \to F$  weakly. Suppose not. Then, there exists a subsequence  $F_{n'_k}$  of  $F_n$  and a continuous and bounded  $f : \mathbb{R} \to \mathbb{R}$  such that, for some  $\varepsilon > 0$ ,

$$\left| \int_{\mathbb{R}} f dF_{n'_k} - \int_{\mathbb{R}} f dF \right| \ge \varepsilon, \quad k = 1, 2, \dots$$
(2.55)

But  $F_{n'_k}$  is tight (being a subsequence of  $F_n$ ), so that again by Helly's theorem there is a subsequence of indeces  $\{n''_k\} \subseteq \{n'_k\}$  such that  $F_{n''_k}$  converges weakly to some distribution Q. Certainly,  $Q \neq F$ , since, otherwise, (2.55) would be violated. Now,  $\varphi_{n_k}(t) \rightarrow \varphi_F(t)$ , the characteristic function of F, and  $\varphi_{n''_k}(t) \rightarrow \varphi_Q(t)$ , where  $\varphi_F(t) \neq \varphi_Q(t)$  since  $F \neq Q$ . These two convergences are contradictory since, by hypothesis, the whole  $\varphi_n(t)$  sequence converges to the same limiting function  $\varphi(t)$ . Thus,  $F_n \rightarrow F$ weakly and this completes the proof.  $\Box$ 

**EXAMPLE 2.12** It is possible that  $\varphi_n(t) \to \varphi(t)$ ,  $\forall t \in \mathbb{R}$ , but:  $\varphi(t)$  is not continuous at t = 0;  $\varphi(t)$  is not a characteristic function; and  $F_n$  is not weakly convergent to a distribution function F. The reader can gain insight in this fact by considering the distribution function  $F_n$  associated with the uniform distribution on [-n,n], whose characteristic function is

$$\varphi_n(t) = \begin{cases} 1, & t = 0\\ \frac{1}{nt}\sin(tn), & t \neq 0. \end{cases}$$
(2.56)

*Here*,  $\varphi_n(t) \rightarrow \varphi(t)$  *with* 

$$\varphi(t) = \begin{cases} 1, & t = 0\\ 0, & t \neq 0, \end{cases}$$
(2.57)

but  $\varphi(t)$  is not a characteristic function (as it can be easily shown, a characteristic function is always a continuous function) and  $F_n$  does not converge weakly to any F.

Thus,  $\varphi_n(t) \rightarrow \varphi(t)$  is not sufficient to conclude that  $F_n$  converges weakly to some F. However, part (b) of Theorem 2.11 tells us that, whenever convergence fails,  $\varphi(t)$  has to be discontinuous in 0.

# 2.4 Gaussian random variables

A n-dimensional random variable is "Gaussian" (or "normal") if it has density function

$$p(x) = \frac{1}{(2\pi)^{n/2} |V|^{1/2}} e^{-\frac{1}{2}(x-m)^T V^{-1}(x-m)},$$
(2.58)

where  $x = [x_1 \ x_2 \ \cdots \ x_n]^T \in \mathbb{R}^n$ ,  $m \in \mathbb{R}^n$ ,  $V \in \mathbb{R}^{n \times n}$  is symmetric and positive definite (we write  $V \succ 0$ ), and  $|\cdot|$  indicates determinant.

It can be computed that

$$m = \mathbb{E}[v]; \tag{2.59}$$

$$V = \mathbb{V}ar(v). \tag{2.60}$$

Thus, the density function of a Gaussian random variable is fully described by its mean and its variance.

Gaussian random variables have notable properties, as listed below.

(i) If an *n*-dimensional random variable v is Gaussian, then, given a matrix  $A \in \mathbb{R}^{m \times n}$  such that  $AA^T \succ 0$ , Av is Gaussian too.

This can be proven by a direct computation, which we here omit. Condition  $AA^T \succ 0$  prevents Av from concentrating in a subspace of  $\mathbb{R}^m$ , in which case the density function of Av does not exist. Also, we have:  $\mathbb{E}[Av] = A\mathbb{E}[v] = Am$  and  $\mathbb{V}ar(Av) = \mathbb{E}[(Av - Am)(Av - Am)^T] = A\mathbb{E}[(v - m)(v - m)^T]A^T = AVA^T$ .

## (ii) Uncorrelation implies independence.

Suppose that the variance matrix *V* has the form

$$V = \begin{bmatrix} V_{11} & 0\\ 0 & V_{22} \end{bmatrix}, \tag{2.61}$$

where  $V_{11}$  and  $V_{22}$  are matrices of size  $n_1 \times n_1$  and  $n_2 \times n_2$ , respectively, that is *v* is formed by two uncorrelated components  $v_1 \in \mathbb{R}^{n_1}$  and  $v_2 \in \mathbb{R}^{n_2}$ . Then, by splitting x - m into two components of suitable dimensions, we have

$$p(x) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2}(x-m)^T V^{-1}(x-m)}$$
(2.62)

$$= \frac{1}{(2\pi)^{n_1/2} |V_{11}|^{1/2} (2\pi)^{n_2/2} |V_{22}|^{1/2}} e^{-\frac{1}{2} [(x_1 - m_1)^T V_{11}^{-1} (x_1 - m_1) + (x_2 - m_2)^T V_{22}^{-1} (x_1 - m_1)^2 (2\pi)^3]}$$

$$= p_1(x_1)p_2(x_2), (2.64)$$

where  $p_1(x_1)$  and  $p_2(x_2)$  are the density functions of  $v_1$  and  $v_2$ , so proving that  $v_1$  and  $v_2$  are independent.

The above two properties are certainly one reason of the success of Gaussian variables: because of these properties, many statistical problems find an easier solution within the Gaussian framework. A second reason of success is that Gaussian variables provide a universal paradigm for the description of natural phenomena that involve many stochastic sources, a fact that has a theoretical foundation in the central limit Theorem (see Section 3.4).

Thus far, we have considered Gaussian variables that have a positive definite variance V. Sometimes it is convenient to have at our disposal a more general definition that allows for a positive semidefinite variance as well. In this case, however, a density function does not exist and we have to move to distribution functions or – as we prefer to do – to characteristic functions.

Let us start by observing that the characteristic function of a Gaussian variable with mean *m* and variance  $V \succ 0$  is given by (we omit the lengthy and conceptually uninteresting derivation):

$$\varphi(t_1, t_2, \dots, t_n) := \mathbb{E}[e^{it^T v}] = e^{it^T m - \frac{1}{2}t^T V t}$$
(2.65)

(in fact we have referred here to the characteristic function of a multidimensional random variable, a definition which naturally extends that valid for the 1-dimensional case. Similarly to Theorem 2.10, multidimensional distributions are in a 1-to-1 correspondence with multidimensional characteristic functions). Suppose now that V is only positive semidefinite and consider again expression

$$e^{it^T m - \frac{1}{2}t^T V t}.$$
 (2.66)

(2.66) still identifies a characteristic function (i.e., it equals  $\int_{\mathbb{R}} e^{it^T x} dF(x)$  for some *n*-dimensional probability distribution *F*). To show this, consider

$$e^{it^T m - \frac{1}{2}t^T (V + \frac{1}{n}I)t},$$
(2.67)

where *I* is the identity matrix. Since  $V + \frac{1}{n}I > 0$ , (2.67) is the characteristic function  $\varphi_n(t)$  of a Gaussian distribution  $G(m, V + \frac{1}{n}I)$ . When we let  $n \to \infty$ , (2.67)  $\to$  (2.66), and the limit  $\varphi(t) = (2.66)$  is continuous at t = 0. Then, by appealing to Theorem 2.11 (actually, to an extension of this theorem to multidimensional distribution functions), we conclude that (2.67) is indeed the characteristic function of some distribution function. This justifies the following definition.

#### **DEFINITION 2.13 (Gaussian random variable)**

A n-dimensional random variable is "Gaussian" (or "normal") if it has characteristic function
$$\varphi(t_1, t_2, \dots, t_n) := E[e^{it^T v}] = e^{it^T m - \frac{1}{2}t^T V t}, \qquad (2.68)$$

where  $t = [t_1 \ t_2 \ \cdots \ t_n]^T, m \in \mathbb{R}^n, \ 0 \preceq V \in \mathbb{R}^{n \times n}$ .

# 2.5 Computing the density induced by a function

Sometimes, it is necessary to compute the density function of a random variable obtained by applying a function f to another random variable whose density function is known. The following theorem provides an answer to this problem.

**THEOREM 2.14 (density function of f**(**v**)) Consider a n-dimensional random variable v with density function  $\mathbb{P}$ . Given a function  $f : \mathbb{R}^n \to \mathbb{R}^n$  such that: i) f is 1-to-1; and ii)  $g = f^{-1}$  is everywhere differentiable, then v' = f(v) is a n-dimensional random variable and it has a density function given by the relation

$$p'(y) = p(g(y)) \cdot |J_g|,$$
 (2.69)

where  $J_g$  is the Jacobian of g, namely

$$J_{g} = \det \begin{bmatrix} \frac{\partial g_{1}}{\partial y_{1}} & \cdots & \frac{\partial g_{1}}{\partial y_{n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_{n}}{\partial y_{1}} & \cdots & \frac{\partial g_{n}}{\partial y_{n}} \end{bmatrix}$$
(2.70)

(subscript denotes component) and  $|\cdot|$  is absolute value.

PROOF. The proof uses results on absolute continuous functions which we take here for granted. Moreover, we only consider the 1-dimensional case since the multi-dimensional case is conceptually similar but notationally more complicated.

For n = 1, the thesis writes

$$p'(y) = p(g(y)) \left| \frac{dg}{dy} \right|.$$
(2.71)

We first establish that v' is a random variable. Start by observing that f, being the inverse of a differentiable and therefore continuous function g, is continuous, so that, by definition of continuity, the inverse image through f of an open set is open, and therefore in  $\mathscr{B}(\mathbb{R})$ . Now, it is not difficult to see that the collection of all sets whose

inverse image through f is in  $\mathscr{B}(\mathbb{R})$  is a  $\sigma$ -algebra. Since Borel sets are the smallest  $\sigma$ -algebra containing the open sets, we conclude that the inverse image of any Borel set is in  $\mathscr{B}(\mathbb{R})$ , i.e. f is  $\mathscr{B}(\mathbb{R})/\mathscr{B}(\mathbb{R})$ -measurable. Hence, v' = f(v) is a random variable in view of Theorem 1.5.

We turn now to prove the validity of (2.71).

Since g is 1-to-1, it is either increasing or decreasing. Suppose g is increasing (the decreasing case goes through similarly). Fix an interval [-M, M], where M is an integer. Then,

$$g(x) - g(-M) = \int_{-M}^{x} \frac{dg}{dy} dy \quad for -M \le x \le M$$
(2.72)

(this follows from Lemma 7.25 and Theorem 7.18 in [6] that prove that a g increasing and everywhere differentiable over [-M, M] is absolutely continuous over the same interval and from the fact that, for an absolutely continuous function g,  $\frac{dg}{dy}$  is measurable and (2.72) holds - Theorem 7.20 in [6]).

Now, put  $\mu(B) = \lambda(g(B)), B \in \mathscr{B}[-M, M]$  ( $\lambda$  is Lebesgue measure). The  $\sigma$ -additivity of  $\lambda$  implies the  $\sigma$ -additivity of  $\mu$ , so that  $\mu$  is a measure on  $\mathscr{B}[-M, M]$ . g(x) - g(-M) is its distribution and, by virtue of (2.72),  $\frac{dg}{dy}$  is its density function (the notions of distribution and density functions used here are the same as in Definitions 2.4 and 2.5 expect for the scaling factor  $\mu[-M, M]$ .) Thus,

$$\lambda(g(B)) = \mu(B) = \int_{[-M,M]} \mathbf{1}(B) d\mu = \int_{[-M,M]} \mathbf{1}(B) \frac{dg}{dy} dy = \int_B \frac{dg}{dy} dy, \quad \forall B \in \mathscr{B}[-M,M],$$
(2.73)

where  $\mathbf{1}(\cdot)$  is the indicator function and the third "=" is justified in view of the comment that follows Definition 2.5.

Turn now to consider the density  $\mathbb{P}$ . Assume first that  $p = 1(A)/\lambda(A)$ , where A is a Borel set with  $\lambda(A) > 0$ . Then

$$\int_{-M}^{x} p(g(y)) \frac{dg}{dy} dy = \frac{1}{\lambda(A)} \int_{[-M,x] \cap f(A)} \frac{dg}{dy} dy \qquad (2.74)$$

$$= \frac{1}{\lambda(A)} \lambda\left(g([-M,x] \cap A)\right) \quad (using \ (2.73)) \qquad (2.75)$$

$$= P\{g(-M) \le v \le g(x)\}$$
(2.76)

$$= P\{-M \le v' \le x\}.$$
(2.77)

Let  $M \to \infty$  to conclude that

$$\int_{-\infty}^{x} p(g(y)) \frac{dg}{dy} dy = P\left\{v' \le x\right\},\tag{2.78}$$

from which we see that  $p(g(y))\frac{dg}{dy}$  is the density of v', that is, (2.71) is established.

In the case of a generic density  $\mathbb{P}$ , to arrive to equation (2.78) one has to first extend the derivation in (2.77) to simple functions, and then pass to the limit by the monotone convergence Theorem 3.8.

Theorem 2.14 can also be applied to functions  $f : \mathbb{R}^n \to \mathbb{R}^m$ , with m < n, provided that  $\mathbb{R}^m$  can be augmented with dummy variables such that the augmented transformation becomes invertible. This is illustrated by an example.

**EXAMPLE 2.15** Given a bi-dimensional v with density function  $\mathbb{P}$ , suppose we want to compute the density function of  $\eta$  defined as the sum of the two components of v:  $\eta = v_1 + v_2$ . Since the transformation  $v \to \eta$  is from  $\mathbb{R}^2$  to  $\mathbb{R}^1$ , Theorem 2.14 cannot be directly applied. However, introducing the dummy variable  $\xi = v_2$  and letting  $v' = [\eta \ \xi]^T$ , we have

$$v' = Av, \tag{2.79}$$

where  $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ , and (2.79) is an invertible transformation. Now,  $A^{-1} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}$ , so that  $|J_g| = \left| \det \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \right| = 1$ , and Theorem 2.14 gives:  $p'(y_1, y_2) = p(y_1 - y_2, y_2).$  (2.80)

*The density function of*  $\eta$  *can then be recovered by integration:* 

$$p_{\eta}(y_1) = \int_{\mathbb{R}} p(y_1 - y_2, y_2) dy_2.$$
(2.81)

The reader is invited to complete this example by computing  $p_{\eta}$  when v is uniformly distributed in  $[0,1]^2$ .

2.5 Computing the density induced by a function

# **Chapter 3**

# **STOCHASTIC CONVERGENCE**

# **3.1** Probabilistic notions of convergence

We introduce a number of probabilistic notions of convergence of a sequence of random variables  $v_n$  to a limit random variable v and relate each one to the others.

**DEFINITION 3.1** (stochastic convergence) Given a sequence of random variables  $v_n$  and an additional random variable v defined on a probability space  $(\Omega, \mathscr{F}, \mathbb{P})$ , we say that  $v_n \to v$ 

(a) uniformly, if  $\sup_{\omega \in \Omega} |v_n(\omega) - v(\omega)| \to 0$ ; (b) surely, if  $v_n(\omega) - v(\omega) \to 0$ ,  $\forall \omega \in \Omega$ ; (c) almost surely, if  $\mathbb{P}\{\omega \text{ such that } v_n(\omega) - v(\omega) \to 0\} = 1$  (when we want to emphasize probability  $\mathbb{P}$  we write  $\mathbb{P}$ -almost surely.) Another expression equivalent to "almost surely" is "with probability 1"; (d) in  $\mathbb{L}^2$ , if  $\mathbb{E}[(v_n - v)^2] \to 0$ ; (e) in  $\mathbb{L}^1$ , if  $\mathbb{E}[|v_n - v|] \to 0$ ; (f) in probability, if  $\forall \varepsilon > 0$ ,  $\mathbb{P}\{\omega \text{ such that } |v_n(\omega) - v(\omega)| \ge \varepsilon\} \to 0$ ; (g) weakly, if for any continuous and bounded function  $f : \mathbb{R} \to \mathbb{R}$ , we have  $\mathbb{E}[f(v_n)] \to \mathbb{E}[f(v)]$ . " $v_n \to v$  weakly" is also expressed as " $v_n \to v$  in distribution". When  $v_n$  converges in distribution to a variable with distribution F, we also write  $v_n \sim AsF$ .  $\Box$ 

Definitions (a)-(f) are concerned with the behavior of  $v_n - v$  and require that this difference goes to zero as  $n \to \infty$  in different ways as specified by the different definitions. Thus, for instance,  $v_n$  tends to v almost surely if  $v_n - v$  tends to zero almost surely. In contrast, the fact that  $v_n \to v$  weakly in no way implies that  $v_n - v \to 0$ . To understand this, suppose e.g. that  $v_n = \xi$ , n = 1, 2, ..., where  $\xi$  is a fixed random variable different from v but sharing with v the same distribution. Then clearly  $\mathbb{E}[f(v_n)] = \mathbb{E}[f(\xi)] = \mathbb{E}[f(v)], \forall n$ , so that  $v_n \to v$  weakly, but  $v_n - v = \xi - v$  does not converge to zero.

Weak convergence is in fact a property of the distribution of the random variables. Indeed,  $\mathbb{E}[f(v_n)] \to \mathbb{E}[f(v)]$  can be rewritten as  $\int_{\mathbb{R}} f dF_n \to \int_{\mathbb{R}} f dF$  (where  $F_n$  is the distribution function of  $v_n$  and F that of v) and we see that weak convergence means that  $F_n$  approaches F. Weak convergence is discussed in detail in Section 3.5.

The different notions of convergence are related to each other by the following theorem.





Figure 3.1: Implications among various notions of stochastic convergence.

**PROOF.** Implications  $(a) \Rightarrow (b) \Rightarrow (c)$  and  $(a) \Rightarrow (d)$  are obvious.

Let 
$$\xi_n := v_n - v$$
.

 $(d) \Rightarrow (e)$ ] By Schwarz inequality 4.7 applied to  $\mathbb{L}^2$  (see Example 4.4):  $\mathbb{E}[|\xi_n|] = \mathbb{E}[1 \cdot |\xi_n|] \le (\mathbb{E}[1^2])^{1/2} (\mathbb{E}[\xi_n^2])^{1/2} = (\mathbb{E}[\xi_n^2])^{1/2}$ , showing that  $\mathbb{E}[\xi_n^2] \to 0$  implies  $\mathbb{E}[|\xi_n|] \to 0$ .  $(c) \Rightarrow (f)$ ] Let  $A_j^{\varepsilon} := \{\omega \text{ such that } |\xi_n| < \varepsilon, \forall n \ge j\}$ .  $A_j^{\varepsilon}$  is increasing with j and  $\bigcup_j A_j^{\varepsilon} =: A^{\varepsilon}$  is the set where the tail of  $|\xi_n|$  is below  $\varepsilon$ . Since  $\xi_n \to 0$  almost surely,  $\mathbb{P}(A^{\varepsilon}) = 1$ , from which  $\mathbb{P}(A_j^{\varepsilon}) \to 1$  as  $j \to \infty$ . Now, since  $\{\omega \text{ such that } |\xi_j| \ge \varepsilon\} \subseteq \Omega - A_j^{\varepsilon}$ , we obtain that  $\mathbb{P}\{\omega \text{ such that } |\xi_j| \ge \varepsilon\} \to 0$ .

 $(e) \Rightarrow (f)$ ] From (e),  $\varepsilon \mathbb{P}\{\omega \text{ such that } |\xi_n| \ge \varepsilon\} \le \mathbb{E}[|\xi_n|] \to 0 \text{ and (f) follows.}$  $(f) \Rightarrow (g)$ ] Given  $\varepsilon_1, \varepsilon_2 > 0$ , fix M and  $\varepsilon$  such that  $\mathbb{P}\{\omega \text{ such that } |v| \ge M\} \le \varepsilon_1$ and  $|f(x) - f(y)| \le \varepsilon_2$  for |y| < M and  $|x - y| < \varepsilon$  (such  $\varepsilon$  exists since a continuous function is uniformly continuous on a bounded set). Then,

$$|\mathbb{E}[f(v_n)] - \mathbb{E}[f(v)]| \tag{3.1}$$

$$\leq \mathbb{E}\left[\left|f(v_n) - f(v)\right|\right] \tag{3.2}$$

$$\leq (2\max_{x} |f(x)|) [\varepsilon_1 + \mathbb{P}\{\omega \text{ such that } |v_n - v| \ge \varepsilon\}] + \varepsilon_2.$$
(3.3)

Since  $\varepsilon_1$  and  $\varepsilon_2$  are arbitrarily small, *f* is bounded, and  $\mathbb{P}\{\omega \text{ such that } |v_n - v| \ge \varepsilon\} \to 0$  by assumption, (g) follows.  $\Box$ 

No other implications than the ones stated in Theorem 3.2 hold true. In particular, almost sure convergence does not imply and is not implied by  $\mathbb{L}^2$ -convergence, as the following example shows.

**EXAMPLE 3.3** Consider the following sequence of random variables defined on the probability space  $([0,1], \mathscr{B}[0,1], \lambda)$ :

$$v_n = \begin{cases} \sqrt{n}, & on [0, 1/n] \\ 0, & otherwise. \end{cases}$$
(3.4)

Letting v := 0, clearly  $v_n \to v$  almost surely, but  $\mathbb{E}[(v_n - v)^2] = \frac{1}{n}n = 1 \not\to 0$ , so that almost sure convergence does not imply  $\mathbb{L}^2$ -convergence.

Conversely, consider the sequence  $v_1^1, v_2^1, v_2^2, v_3^1, v_3^2, v_3^3, \dots$  with

$$v_n^k = \begin{cases} 1, & on\left[\frac{k-1}{n}, \frac{k}{n}\right] \\ 0, & otherwise. \end{cases}$$
(3.5)

This sequence is  $\mathbb{L}^2$ -convergent to zero, but, for every  $\omega \in [0,1]$ , the sequence keeps oscillating between 0 and 1 so that it does not converge for any  $\omega$ .

The reason why we had  $\mathbb{L}^2$ -convergence but not almost sure convergence in the latter example was that the intervals where  $v_n^k = 1$  in (3.5) had two properties: i) their size shrinks (so that  $\mathbb{L}^2$  convergence to zero takes place); and ii) each point in [0,1] falls infinitely many times in the intervals (and, thus, almost sure convergence fails). Fact ii) is possible because the sum of the interval lengths  $1, \frac{1}{2}, \frac{1}{2}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \dots$  where  $v_n^k = 1$  is diverging. The following theorem, which shows a converse result that when certain sums are finite then almost sure convergence takes place, is important to assess almost sure convergence in various contexts.

**THEOREM 3.4** Let  $v_n$ , n = 1, 2, ..., and v be random variables. Suppose that for any  $\varepsilon > 0$ ,

$$\sum_{n=1}^{\infty} \mathbb{P}\{\omega \text{ such that } |v_n - v| \ge \varepsilon\} < \infty,$$
(3.6)

then  $v_n \rightarrow v$  almost surely.

PROOF. Letting

$$A_j^k := \{ \omega \text{ such that } |v_n - v| \ge \frac{1}{k}, \text{ for some } n \ge j \},$$
(3.7)

we have that  $\{\omega \text{ such that } v_n - v \not\to 0\} = \bigcup_{k=1}^{\infty} \bigcap_{j=1}^{\infty} A_j^k$ . Thus,

$$\mathbb{P}\{\boldsymbol{\omega} \text{ such that } \boldsymbol{v}_n - \boldsymbol{v} \neq 0\} \leq \sum_{k=1}^{\infty} \lim_{j \to \infty} \sum_{n=j}^{\infty} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_n - \boldsymbol{v}| \geq \frac{1}{k}\}.$$
(3.8)

Since (3.6) holds, each single term  $\lim_{j\to\infty}\sum_{n=j}^{\infty}\mathbb{P}\{\omega \text{ such that } |v_n - v| \ge \frac{1}{k}\}$  is zero and the right-hand side of the previous inequality is null, so proving that  $v_n \to v$  almost surely.

### **3.2** Measurability of the limit of random variables

The next result relates the measurability of a sequence of random variables to that of its limit.

**THEOREM 3.5** Let  $v_n$  be a sequence of random variables on  $(\Omega, \mathscr{F}, \mathbb{P})$  (so that each  $v_n$  is  $\mathscr{F}$ -measurable) and let v be an additional variable that is not required to be  $\mathscr{F}$ -measurable by assumption. If  $v_n(\omega) \to v(\omega)$ ,  $\forall \omega \in \Omega$ , then v is  $\mathscr{F}$ -measurable.

**PROOF.** For any  $a, b \in \mathbb{R}$ , we have

$$\{\boldsymbol{\omega} \text{ such that } \boldsymbol{v} \in (a,b)\} = \bigcup_{p=1}^{\infty} \bigcap_{n>p} \{\boldsymbol{\omega} \text{ such that } \boldsymbol{v}_n \in (a,b)\}.$$
(3.9)

Since  $\{\omega \text{ such that } v_n \in (a,b)\} \in \mathscr{F}$  and a  $\sigma$ -algebra is closed under countable intersection and union, we have that  $\{\omega \text{ such that } v \in (a,b)\} \in \mathscr{F}$ , from which the measurability of *v* follows by applying the test of measurability 2.3.

#### The almost sure limit of $v_n$

The set  $\{v_n \rightarrow\}$  where the limit of a sequence of random variables exists finite is always measurable, i.e., it belongs to  $\mathscr{F}$ .

To see that this is the case, le us consider the complementary set where  $v_n$  does not converge to a finite limit, i.e., where  $v_n$  diverges or it oscillates. The set where  $v_n \to \infty$  can be written as  $\bigcap_{M=1}^{\infty} \bigcup_{p=1}^{\infty} \bigcap_{n \ge p} \{\omega \text{ such that } v_n \ge M\}$ , where the curly bracket is the

set of the  $\omega$ 's where, for any fixed M, the tail of  $v_n$  from a certain p onward is above M. Since  $\{\omega \text{ such that } v_n \ge M\}$  is measurable and a  $\sigma$ -algebra is closed under countable intersection and union, we have that  $\{\omega \text{ such that } v_n \to \infty\}$  is measurable. Similarly,  $\{\omega \text{ such that } v_n \to -\infty\}$  is measurable. Consider then a  $\omega$  where  $v_n(\omega)$  oscillates. Then, there exist two rational numbers  $\alpha$  and  $\beta$  with  $\alpha < \beta$  such that  $v_n(\omega)$  is below  $\alpha$  infinitely many times and above  $\beta$  infinitely many times. The set of  $\omega$ 's where  $v_n(\omega) < \alpha$  infinitely many times can be written as  $\bigcap_{p=1}^{\infty} \bigcup_{n \ge p} \{\omega \text{ such that } v_n < \alpha\}$ and it is measurable. Likewise is measurable the set where  $v_n(\omega) > \beta$  infinitely many times, so that we obtain the measurability of the set where  $v_n(\omega)$  oscillates between  $\alpha$  and  $\beta$  infinitely many times, which is the intersection of the two previous sets. The set where  $v_n$  oscillates is obtained as union over all rationals  $\alpha$  and  $\beta$  and is therefore measurable. In conclusion, the set where  $v_n$  does not converge to a finite limit is the union of measurable sets and is measurable, and so its complementary set  $\{v_n \to \}$  is also measurable.

Suppose now that the measurable set  $\{v_n \rightarrow\}$  has probability 1:  $\mathbb{P}\{v_n \rightarrow\} = 1$ , that is, the sequence  $v_n$  admits almost surely a limit. If we let

$$\bar{v}_n = \begin{cases} v_n, & \text{where } v_n \text{ converges to a finite value} \\ 0, & \text{otherwise}, \end{cases}$$
(3.10)

this  $\bar{v}_n$  is measurable and it converges for any  $\omega \in \Omega$  to

$$v = \begin{cases} \lim_{n \to \infty} v_n, & \text{where } v_n \text{ converges to a finite value} \\ 0, & \text{otherwise,} \end{cases}$$
(3.11)

and, by an application of Theorem 3.5 we see that such a v is a random variable. v is an "almost sure limit of  $v_n$ ". Notice also that any other random variable obtained as  $v + \eta$ , with  $\eta = 0$  almost surely, is also an almost sure limit of  $v_n$ . Since two almost sure limits only differ on a zero probability set, to many purposes specifying which limit one is considering is immaterial, and it is customary to speak of "almost sure limit of  $v_n$ ", where it is meant that one refers to anyone among the random variables that are almost sure limits of  $v_n$ .

For easy reference, we summarize the discussion in the following theorem.

**THEOREM 3.6** Consider a sequence of random variables  $v_n$ . Then,

$$v = \begin{cases} \lim_{n \to \infty} v_n, & \text{where } v_n \text{ converges to a finite value} \\ 0, & \text{otherwise,} \end{cases}$$
(3.12)

is a random variable (i.e. it belongs to  $\mathscr{F}$ ). The set  $\{v_n \rightarrow\}$  where  $v_n$  converges to a finite value is measurable. If  $\mathbb{P}\{v_n \rightarrow\} = 1$ , then v is called the almost sure limit of  $v_n$ . Any other random variable  $v + \eta$  obtained by adding to v a random variable  $\eta$  with  $\eta = 0$  almost surely is also called an almost sure limit of  $v_n$ .

Measurability with respect to a sub  $\sigma$ -algebra  $\mathscr{G} \subseteq \mathscr{F}$ 

**THEOREM 3.7** Let  $v_n$  be a sequence of random variables on  $(\Omega, \mathscr{F}, \mathbb{P})$  that are  $\mathscr{G}$ -measurable for some  $\sigma$ -algebra  $\mathscr{G} \subseteq \mathscr{F}$  and let v be an additional random variable. If  $v_n \to v$  in probability, then v need not be  $\mathscr{G}$ -measurable, but there exists a  $\mathscr{G}$ -measurable  $\bar{v}$  such that  $\mathbb{P}\{v \neq \bar{v}\} = 0$ .

**PROOF.** Fix a sequence of real numbers  $\varepsilon_k \downarrow 0$  (i.e.,  $\varepsilon_k$  is decreasing and tends to zero) and extract from  $v_n$  a subsequence  $v_{n_k}$  such that

$$\sum_{k=1}^{\infty} \mathbb{P}\{\omega \text{ such that } |v_{n_k} - v| \ge \varepsilon_k\} < \infty,$$
(3.13)

(such a sequence exists since  $v_n \to v$  in probability). Equation (3.13) implies that  $v_{n_k} \to v$  almost surely, as it can be proven by applying Theorem 3.4. Indeed, given  $\varepsilon > 0$ , let  $\bar{k}$  be such that  $\varepsilon_{\bar{k}} \le \varepsilon$  and write

$$\sum_{k=1}^{\infty} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}\}$$

$$= \sum_{k=1}^{\bar{k}-1} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}\} + \sum_{k=\bar{k}}^{\infty} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}\}$$
(3.14)
$$\leq \sum_{k=1}^{\bar{k}-1} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}\} + \sum_{k=\bar{k}}^{\infty} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}\}$$
(3.15)
$$\leq \sum_{k=1}^{\bar{k}-1} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}\} + \sum_{k=\bar{k}}^{\infty} \mathbb{P}\{\boldsymbol{\omega} \text{ such that } |\boldsymbol{v}_{n_{k}} - \boldsymbol{v}| \ge \boldsymbol{\varepsilon}_{k}\}$$
(3.16)
$$< \infty,$$
(3.17)

so that the assumption (3.6) of Theorem 3.4 is satisfied.

Now,  $v_{n_k}$  is  $\mathscr{G}$ -measurable and therefore we can see  $v_{n_k}$  as a sequence of random variables on the probability space  $(\Omega, \mathscr{G}, \mathbb{P})$ . Moreover,  $v_{n_k}$  is almost surely convergent to v, so that  $\mathbb{P}\{v_{n_k} \to\} = 1$ . The almost sure limit  $\bar{v}$  of  $v_{n_k}$  is also  $\mathscr{G}$ -measurable and it coincides almost surely with v.

The reader may have noticed that the reason why we have to introduce  $\bar{v}$  is that v can possibly exhibit some "strange" behavior where  $v_{n_k} \not\rightarrow v$  so that v is not  $\mathscr{G}$ -measurable.

## **3.3** Limit under the sign of expectation

Suppose that  $v_n \to v$  almost surely. Under what conditions is it true that  $\mathbb{E}[v_n] \to \mathbb{E}[v]$ ? The following theorems provide an answer.

**THEOREM 3.8 (monotone convergence)** Let  $v_n, n = 1, 2, ..., and v$  be random variables such that  $v_n \uparrow v$  almost surely (i.e.  $v_n$  is increasing and tends to v almost surely), and assume that  $v_n \ge z, n = 1, 2, ...,$  for some random variable z with  $\mathbb{E}[z] > -\infty$ . Then,

$$\mathbb{E}[v_n] \uparrow \mathbb{E}[v]. \tag{3.18}$$

**THEOREM 3.9 (dominated convergence)** Let  $v_n, n = 1, 2, ..., and v$  be random variables such that  $v_n \rightarrow v$  almost surely, and assume that  $|v_n| \leq z, n = 1, 2, ...,$  for some random variable z with  $\mathbb{E}[z] < \infty$ . Then,

$$\mathbb{E}[v_n] \to \mathbb{E}[v]. \tag{3.19}$$

A proof of these theorems can be found in any textbook on probability.

In the statements of the theorems, two types of conditions are present:  $v_n$  is required to approach v; and  $v_n$  is bounded by z. The latter condition serves the purpose to limit the importance of the mismatch between  $v_n$  and v on events of small probability. An example clarifies this matter.

**EXAMPLE 3.10 (Example 3.3 continued)** Consider again the  $v_n$ 's in (3.4). Clearly,  $v_n^2 \rightarrow 0$  almost surely, but  $\mathbb{E}[v_n^2] = 1 \not\rightarrow \mathbb{E}[0] = 0$ . Here, no dominating z exists with  $\mathbb{E}[z] < \infty$ , so that the conditions of Theorem 3.9 are violated.

# **3.4** Convergence results for independent random variables

We commence by proving probabilistic inequalities. Besides being useful to prove convergence results, these inequalities are of interest in their own right.

#### **MARKOV'S INEQUALITY 3.11**

For any nonnegative random variable *v* and real number  $\varepsilon > 0$ ,

$$\mathbb{P}\{v \ge \varepsilon\} \le \frac{\mathbb{E}[v]}{\varepsilon}.$$
(3.20)

PROOF. The proof is elementary:  $\mathbb{E}[v] = \int_{\Omega} v d\mathbb{P} \ge \int_{\{v \ge \varepsilon\}} v d\mathbb{P} \ge \varepsilon \mathbb{P}\{v \ge \varepsilon\}.$   $\Box$ 

An application of Markov's inequality gives

#### **CHEBYSHEV'S INEQUALITY 3.12**

For any  $\varepsilon > 0$ ,

$$\mathbb{P}\{|v| \ge \varepsilon\} \le \frac{\mathbb{E}[v^2]}{\varepsilon^2}.$$
(3.21)

PROOF.

$$\mathbb{P}\{|v| \ge \varepsilon\} = \mathbb{P}(v^2 \ge \varepsilon^2) \tag{3.22}$$

$$\leq \frac{\mathbb{E}[\nu^2]}{\varepsilon^2} \quad (use \ (3.20)). \tag{3.23}$$

1 1		- E
	- 1	
	_	- L

In Markov's inequality, the idea is to lowerbound  $\mathbb{E}[v]$  by squeezing the tail of v to the boundary value  $\varepsilon$ . Thus, the bound is tight only when the tail rapidly vanishes after  $\varepsilon$ . A similar observation applies to Chebyshev's inequality. Better bounds can be found by "redressing" the random variable distribution through some transformation before Markov's inequality is applied. One such example is given by the following inequality due to Chernoff. In this inequality, s is a free parameter that can be used to tune the distribution obtained after transformation and an example of use of s is found in the proof of Hoeffding's inequality (Theorem 3.15.)

#### **CHERNOFF'S INEQUALITY 3.13**

For any s > 0 and  $\varepsilon > 0$ ,

$$\mathbb{P}\{v \ge \varepsilon\} \le \frac{\mathbb{E}[e^{sv}]}{e^{s\varepsilon}}.$$
(3.24)

PROOF.

$$\mathbb{P}\{v \ge \varepsilon\} = \mathbb{P}\{e^{sv} \ge e^{s\varepsilon}\}$$
(3.25)

$$\leq \frac{\mathbb{E}[e^{sv}]}{e^{s\varepsilon}}. \quad (use \ (3.20)). \tag{3.26}$$

#### **Concentration inequalities**

Consider a sequence of independent random variables  $v_k, k = 1, 2, ...$  Concentration inequalities study how a function  $f(v_1, v_2, ..., v_n)$  of the first *n* variables in the sequence concentrates around its expected value  $\mathbb{E}[f(v_1, v_2, ..., v_n)]$ .

Here, we are mainly concerned with the deviation of the normalized sum of random variables (empirical mean) from its mean, that is, our interest is on function  $f(v_1, v_2, ..., v_n) = \frac{1}{n} \sum_{k=1}^{n} v_k$  and we study the behavior of

$$M_n := \frac{1}{n} \sum_{k=1}^n v_k - \mathbb{E}\left[\frac{1}{n} \sum_{k=1}^n v_k\right].$$
 (3.27)

A first bound is obtained by means of Chebyshev's inequality:

$$\mathbb{P}\{|M_n| \ge \varepsilon\} \le \frac{\mathbb{E}[M_n^2]}{\varepsilon^2} = \frac{\frac{1}{n^2} \sum_{k=1}^n \mathbb{V}ar(v_k)}{\varepsilon^2}.$$
(3.28)

**EXAMPLE 3.14** For an independent and identically distributed sequence of Bernoulli random variables (i.e.  $\mathbb{P}\{v_k = 1\} = 1 - \mathbb{P}\{v_k = 0\} = p$ ), from (??) we have

$$\mathbb{P}\left\{ \left| \frac{1}{n} \sum_{k=1}^{n} v_k - p \right| \ge \varepsilon \right\} \le \frac{p(1-p)}{n\varepsilon^2}.$$
(3.29)

Do we expect that bound (3.28) is tight? (remember that Chebyshev's inequality is tight when the distribution tail vanishes rapidly after  $\varepsilon$ ). Applying the central limit Theorem 3.20 leads to the conclusion that, under mild assumptions, the distribution of  $M_n$  tends weakly to a Gaussian, a long-tailed distribution. For example, in the case of the Bernulli sequence of Example 3.14, letting  $\Phi(x) = \int_{-\infty}^{x} (2\pi)^{-1/2} e^{-r^2/2} dr$ , the central limit theorem states that

$$\mathbb{P}\left\{\sqrt{\frac{n}{p(1-p)}}\left(\frac{1}{n}\sum_{k=1}^{n}v_k-p\right)\geq x\right\}\to 1-\Phi(x)\leq \frac{e^{-x^2/2}}{x\sqrt{2\pi}},\qquad(3.30)$$

from which we would expect something like

$$\mathbb{P}\left\{ \left| \frac{1}{n} \sum_{k=1}^{n} v_k - p \right| \ge \varepsilon \right\} \sim e^{-\frac{n\varepsilon^2}{2p(1-p)}},$$
(3.31)

that is the probability decays exponentially fast with n. The gap between inverse-linear (as in (3.29)) and exponential convergence in n is truly large; this gap can be filled in by resorting to Hoeffding's inequality.

**THEOREM 3.15 (Hoeffding's inequality)** Let  $v_k, k = 1, 2, ...,$  be independent bounded random variables taking value in  $[\alpha_k, \beta_k]$  and let  $S_n$  be defined as in (3.27). Then, for any  $\varepsilon > 0$ ,

$$\mathbb{P}\{M_n \ge \varepsilon\} \le e^{-\frac{2n^2\varepsilon^2}{\sum_{k=1}^n (\beta_k - \alpha_k)^2}};$$
(3.32)

and

$$\mathbb{P}\{M_n \le -\varepsilon\} \le e^{-\frac{2n^2\varepsilon^2}{\sum_{k=1}^n (\beta_k - \alpha_k)^2}}.$$
(3.33)

PROOF. For ease of notation, we assume  $[\alpha_k, \beta_k] = [0, 1]$ , in which case we prove that

$$\mathbb{P}\{M_n \ge \varepsilon\} \le e^{-2n\varepsilon^2}; \tag{3.34}$$

and

$$\mathbb{P}\{M_n \le -\varepsilon\} \le e^{-2n\varepsilon^2}.$$
(3.35)

The extension is easy.

We start by observing that, for any random variable *v* with  $\mathbb{E}[v] = 0$  and  $\alpha \le v \le 1 + \alpha$  and for any h > 0, it holds that

$$\mathbb{E}[e^{hv}] \le e^{\frac{h^2}{8}}.\tag{3.36}$$

In fact, by convexity of the exponential function,  $e^{hv} \le (v-\alpha)e^{(1+\alpha)h} + (1+\alpha-v)e^{\alpha h}$ , so that

$$\mathbb{E}[e^{hv}] \leq \mathbb{E}\left[(v-\alpha)e^{(1+\alpha)h} + (1+\alpha-v)e^{\alpha h}\right]$$
(3.37)

$$= \mathbb{E}\left[-\alpha e^{(1+\alpha)h} + (1+\alpha)e^{\alpha h}\right] \quad (since \mathbb{E}[v] = 0) \quad (3.38)$$

$$= -\alpha e^{(1+\alpha)h} + (1+\alpha)e^{\alpha h}$$
(3.39)

$$= e^{\Gamma(h)}, \tag{3.40}$$

where  $\Gamma(h) = \alpha h + \ln(1 + \alpha - \alpha e^h)$ . The derivative of  $\Gamma(h)$  is  $\Gamma'(h) = \alpha - \alpha/[(1 + \alpha)e^{-h} - \alpha]$ , so that  $\Gamma'(0) = 0$ . Moreover,

$$\Gamma''(h) = \frac{-\alpha(1+\alpha)e^{-h}}{[(1+\alpha)e^{-h}-\alpha]^2}$$
(3.41)

$$= \frac{ab}{[a+b]^2} \quad (where we let a = (1+\alpha)e^{-h}, b = -\alpha) \tag{3.42}$$

$$\leq \frac{1}{4}, \quad \forall h. \tag{3.43}$$

Thus, by Taylor series expansion, for some  $\xi \in [0,h]$ :

$$\Gamma(h) = \Gamma(0) + \Gamma'(0)h + \frac{1}{2}\Gamma''(\xi)h^2 \le \frac{h^2}{8},$$
(3.44)

which, used in (3.40), yields (3.36).

Thanks to (3.36), equation (3.34) is now easily obtained from Chernoff's inequality:

$$\mathbb{P}\{S_n \ge \varepsilon\} \le \frac{\mathbb{E}[e^{sS_n}]}{e^{s\varepsilon}} \quad (use \ Chernoff's \ inequality \ (3.24)) \tag{3.45}$$

$$= \frac{\mathbb{E}\left[e^{s\frac{1}{n}\sum_{k=1}^{n}\left(v_{k}-E\left[v_{k}\right]\right)}\right]}{e^{s\varepsilon}}$$
(3.46)

$$= \frac{\prod_{k=1}^{n} \mathbb{E}[e^{s\frac{1}{n}(v_k - E[v_k])}]}{e^{s\varepsilon}} \quad (by \ the \ independence \ of \ the \ v'_k s) \ (3.47)$$

$$\leq \frac{\prod_{k=1}^{n} e^{\frac{s}{8n^2}}}{e^{s\varepsilon}} \quad (use \ (3.36) \ with \ h = s/n) \tag{3.48}$$

$$\leq e^{-2n\varepsilon^2}$$
 (choose  $s = 4\varepsilon n$ ). (3.49)

Equation (3.35) is obtained similarly.

47

**EXAMPLE 3.16 (Example 3.14 continued)** Using Hoeffding's inequality yields

$$\mathbb{P}\left\{ \left| \frac{1}{n} \sum_{k=1}^{n} v_k - p \right| \ge \varepsilon \right\} \le 2e^{-2n\varepsilon^2}$$
(3.50)

(*compare with* (3.31).)

Hoeffding's inequality deals specifically with empirical means, showing that the empirical mean rapidly concentrates around the true mean value. The reason why this is so is that in the empirical mean each single variable has a moderate influence on the computation of the overall empirical mean value and, moreover, different variables do not cooperate because they are independent.

It is a fact that Hoeffding's inequality can be extended to more general functions provided that each variable has a marginal importance in determining the value of the function, as the following theorem states (for a proof see e.g. [4].)

**THEOREM 3.17 (the bounded difference inequality)** Let  $v_k, k = 1, 2, ..., be$ independent random variables taking value in a set A and assume that  $\forall x_1, ..., x_n \in A$ ,  $x'_k \in A$ , and  $\forall k \in [1, n]$ , the measurable function  $f : \mathbb{R}^n \to \mathbb{R}$  satisfies the condition:

$$|f(x_1,\ldots,x_{k-1},x_k,x_{k+1},\ldots,x_n) - f(x_1,\ldots,x_{k-1},x'_k,x_{k+1},\ldots,x_n)| \le \gamma_k.$$
(3.51)

*Then, for any*  $\varepsilon > 0$ *,* 

$$\mathbb{P}\{f(v_1, v_2, \dots, v_n) - \mathbb{E}[f(v_1, v_2, \dots, v_n)] \ge \varepsilon\} \le e^{-\frac{2\varepsilon^2}{\sum_{k=1}^n \gamma_k^2}};$$
(3.52)

and

$$\mathbb{P}\{f(v_1, v_2, \dots, v_n) - \mathbb{E}[f(v_1, v_2, \dots, v_n)] \le -\varepsilon\} \le e^{-\frac{2\varepsilon^2}{\sum_{k=1}^n \gamma_k^2}};$$
(3.53)

Note that (3.52) and (3.53) reduce to (3.32) and (3.33) when we consider empirical means.

#### Laws of large numbers

The laws of large numbers study the convergence of  $\frac{1}{n}\sum_{k=1}^{n} v_k$  to  $\mathbb{E}\left[\frac{1}{n}\sum_{k=1}^{n} v_k\right]$ . This is probably the most studied problem in probability theory and the literature offers an abundant supply of results under various assumptions and according to different notions of convergence. Here, we only present a standard result and prove it by means of concentration inequalities.

**THEOREM 3.18 (law of large numbers)** Let  $v_k, k = 1, 2, ..., be$  independent random variables with uniformly bounded variance:  $var(v_k) \leq C, \forall k$ . Then,

$$\frac{1}{n}\sum_{k=1}^{n}v_{k} \to \mathbb{E}\left[\frac{1}{n}\sum_{k=1}^{n}v_{k}\right] \quad almost \ surely. \tag{3.54}$$

Before proving the theorem, we would like to note that the result is immediate from Hoeffding's inequality if the variables are uniformly bounded. In fact if  $v_k \in [\alpha, \beta]$ ,  $\forall k$ , then

$$\mathbb{P}\{|M_n| \ge \varepsilon\} \le 2e^{-\frac{2n\varepsilon^2}{(\beta-\alpha)^2}},\tag{3.55}$$

where  $M_n := \frac{1}{n} \sum_{k=1}^n v_k - \mathbb{E} \left[ \frac{1}{n} \sum_{k=1}^n v_k \right]$ , so that almost sure convergence  $M_n \to 0$  almost surely follows from Theorem 3.4. On the other hand, if we only assume the boundedness of the variance of the  $v_k$ 's (as is done in the theorem), Hoeffding's inequality does not apply. On the other hand, resorting to Chebyshev's inequality (3.28) yields

$$\mathbb{P}\{|M_n| \ge \varepsilon\} \le \frac{\mathbb{E}[M_n^2]}{\varepsilon^2} = \frac{C}{n\varepsilon^2},\tag{3.56}$$

which is not enough to prove that  $M_n \to 0$  almost surely by way of Theorem 3.4 since  $\sum_{n=1}^{\infty} \frac{C}{n\epsilon^2} = \infty$ . The proof of Theorem 3.18 given below suggests a way to get around this difficulty.

Instead of directly proving the theorem, we prefer to state the following Lemma 3.19 (from which the theorem immediately follows) because the lemma is more general and useful in other contexts as well.

**LEMMA 3.19** Consider the doubly indexed set of random variables  $S_r^p$  such that  $S_r^p = 0$  for r > p and assume that  $\mathbb{E}[(S_r^p)^2] \le C(p+1-r)$  for some constant C and that, for m < n,  $|S_1^n| \le |S_1^m| + |S_{m+1}^n|$ . Then,

$$\frac{1}{n}S_1^n \to 0 \quad almost \ surely. \tag{3.57}$$

Note that Theorem 3.18 immediately follows from the lemma by the position  $S_r^p := \sum_{k=r}^p (v_k - E[v_k])$ .

PROOF OF THE LEMMA. Given an integer *n*, let *N* be the integer such that  $N^2 \le n < (N+1)^2$  and write:

$$\left|\frac{1}{n}S_{1}^{n}\right| \leq \frac{1}{N^{2}}\left|S_{1}^{N^{2}}\right| + \frac{1}{N^{2}}\left|S_{N^{2}+1}^{n}\right|.$$
(3.58)

The lemma is proven by showing that both terms in the last expression go to zero almost surely.

As for the first term, by Chebyshev's inequality we have

$$\sum_{N=1}^{\infty} \mathbb{P}\left\{\frac{1}{N^2} \left|S_1^{N^2}\right| \ge \varepsilon\right\} \le \sum_{N=1}^{\infty} \frac{\mathbb{V}ar\left(\frac{1}{N^2} \left|S_1^{N^2}\right|\right)}{\varepsilon^2}$$
(3.59)

$$\leq \sum_{N=1}^{\infty} \frac{CN^2}{N^4 \varepsilon^2}$$
(3.60)

$$< \infty,$$
 (3.61)

from which almost sure convergence to zero follows using Theorem 3.4. For the second term in (3.58), again using Chebyshev's inequality, we instead have

$$\sum_{n=1}^{\infty} \mathbb{P}\left\{\frac{1}{N^2} \left|S_{N^2+1}^n\right| \ge \varepsilon\right\} \le \sum_{n=1}^{\infty} \frac{\mathbb{V}ar\left(\frac{1}{N^2} \left|S_{N^2+1}^n\right|\right)}{\varepsilon^2}$$
(3.62)

$$\leq \sum_{n=1}^{\infty} \frac{C(n-N^2)}{N^4 \varepsilon^2}$$
(3.63)

$$\leq \sum_{N=1}^{\infty} \sum_{n=N^2}^{(N+1)^2 - 1} \frac{C(n-N^2)}{N^4 \varepsilon^2}$$
(3.64)

$$\leq \sum_{N=1}^{\infty} ((N+1)^2 - N^2) \frac{C((N+1)^2 - 1 - N^2)}{N^4 \varepsilon^2} (3.65)$$

$$\leq \sum_{N=1}^{\infty} (2N+1) \frac{C2N}{N^4 \varepsilon^2}$$
(3.66)

and, again, Theorem 3.4 can be resorted to to prove almost sure convergence to zero. Hence, the two terms in the right-hand side of (3.58) tends to zero almost surely and

this completes the proof.

#### **Central limit theorems**

The literature on central limit theorems is truly vast. We only provide some fundamental results.

Theorem 3.18 shows that, for independent random variables with boundend variance, the following convergence takes place

$$\frac{1}{n}\sum_{k=1}^{n}v_{k} \to \mathbb{E}\left[\frac{1}{n}\sum_{k=1}^{n}v_{k}\right] \quad almost \ surely, \tag{3.68}$$

or, equivalently,

$$\frac{1}{n}\sum_{k=1}^{n}\left(v_{k}-\mathbb{E}[v_{k}]\right)\to 0 \quad almost \ surely. \tag{3.69}$$

A question that arises naturally is as to how fast convergence to zero takes place. This question is answered by the following central limit theorem.

**NOTE:** In condition (3.70) of the theorem, the integral has to be intended as follows. Consider the total variation H(x) of function  $f := F_k - \Phi_k$ , namely H(x) := $\sup \sum_{i=1}^{N} |f(t_i) - f(t_{i-1})|$  where supremum is taken over all N and over all choices of  $t_i$  such that  $t_0 < t_1 < \cdots < t_n = x$ . H(x) is nondecreasing and tends to a finite constant  $\alpha$  as  $x \to \infty$ . If  $\alpha = 0$  (which only happens if  $F_k = \Phi_k$ ), then the integral is taken with respect to the zero measure and its value is zero. Otherwise,  $H(x)/\alpha$  is a probability distribution function. Integration is with respect to this measure and, to compensate for the division by  $\alpha$ , the integrand is multiplied by this same  $\alpha$  value.

**THEOREM 3.20 (central limit theorem)** Let  $v_k, k = 1, 2, ..., be$  independent random variables with probability distribution function  $F_k$ , zero mean and variance  $\sigma_k^2$ , and let  $V_n^2 := \sum_{k=1}^n \sigma_k^2$ . Assume that, for n large enough,  $V_n^2 > 0$ , that is, at least one variable  $v_k$  has non-zero variance (this only serves the purpose of avoiding division by zero). Moreover, let  $\Phi_k(x) = \int_{-\infty}^x (2\pi)^{-1/2} \sigma_k^{-1} e^{-r^2/2\sigma_k^2} dr$  be the Gaussian distribution function with zero mean and the same variance as  $v_k$ . If

$$\forall \varepsilon > 0, \quad \frac{1}{V_n^2} \sum_{k=1}^n \int_{|x| > \varepsilon V_n} x^2 d|F_k - \Phi_k| \to 0, \quad as \ n \to \infty$$
(3.70)

(see the NOTE before the theorem on how this integral must be interpreted), then,

$$\frac{1}{V_n} \sum_{k=1}^n v_k \sim AsG(0,1).$$
(3.71)

Based on the results in Appendix 2.4, we know that the sum of jointly Gaussian random variables is Gaussian too. In words, Gaussianity is a "closed world": once we are in it, we cannot get out by applying linear operations. Theorem 3.20 tells us more: this world is also "attractive" and the sum of independent variables tends to be Gaussian under general conditions.

Before proving the theorem we make an observation on the theorem assumptions.

#### Observation

Assumption (3.70) in the theorem is implied by each of the following handier sufficient conditions:

1. 
$$\forall \varepsilon > 0$$
,  $\frac{1}{V_n^2} \sum_{k=1}^n \int_{|x| > \varepsilon V_n} x^2 dF_k(x) \to 0$ , as  $n \to \infty$  (Lindeberg's condition);  
2.  $\frac{1}{V_n^{2+\delta}} \sum_{k=1}^n \mathbb{E}[|v_k|^{2+\delta}] \to 0$ , for some  $\delta > 0$  (Lyapunov's condition).

We show that Lyapunov's condition implies Lindeberg's condition which, in turn, implies (3.70).

Lyapunov's condition  $\Rightarrow$  Lindeberg's condition) We have:

$$\mathbb{E}[|v_k|^{2+\delta}] \geq \int_{|x| > \varepsilon V_n} |x|^{2+\delta} dF_k(x)$$
(3.72)

$$\geq \varepsilon^{\delta} V_n^{\delta} \int_{|x| > \varepsilon V_n} x^2 dF_k(x), \qquad (3.73)$$

which gives

$$\frac{1}{V_n^2} \sum_{k=1}^n \int_{|x| > \varepsilon V_n} x^2 dF_k(x) \leq \frac{1}{\varepsilon^{\delta} V_n^{2+\delta}} \sum_{k=1}^n E[|v_k|^{2+\delta}]$$
(3.74)

$$\rightarrow 0$$
, (useLyapunov's condition) (3.75)

that is Lindeberg's condition.

Lindeberg's condition  $\Rightarrow$  (3.70)) Note first that Lindeberg's condition implies

$$\frac{\max_{1 \le k \le n} \sigma_k^2}{V_n^2} \to 0.$$
(3.76)

Indeed,

$$\frac{\max_{1 \le k \le n} \sigma_k^2}{V_n^2} \le \frac{\sum_{k=1}^n \int_{|x| > \varepsilon V_n} x^2 dF_k(x) + \varepsilon^2 V_n^2}{V_n^2}$$
(3.77)

$$\rightarrow \epsilon^2,$$
 (3.78)

where Lindeberg's condition has been applied in computing the limit. Since  $\varepsilon$  is arbitrary, (3.76) follows. Now, letting  $\Phi(x)$  be the Gaussian distribution with zero mean and unitary variance, we have:

$$\frac{1}{V_n^2} \sum_{k=1}^n \int_{|x| > \varepsilon V_n} x^2 d\Phi_k(x) = \frac{1}{V_n^2} \sum_{\substack{k=1,\dots,n \\ \sigma_k \neq 0}}^n \int_{|x| > \varepsilon V_n} x^2 d\Phi_k(x)$$
(3.79)

$$= \frac{1}{V_n^2} \sum_{k=1}^n \int_{|z| > \frac{\varepsilon V_n}{\sigma_k}} \sigma_k^2 z^2 d\Phi(z) \quad (where \ z = x/\sigma_k) 3.80)$$

$$\leq \frac{1}{V_n^2} \sum_{k=1}^n \int_{|z| > \frac{\varepsilon V_n}{\max_{1 \le k \le n} \sigma_k}} \sigma_k^2 z^2 d\Phi(z)$$
(3.81)

$$= \int_{|z| > \frac{\varepsilon_{V_n}}{\max_{1 \le k \le n} \sigma_k}} z^2 d\Phi(z) \cdot \frac{1}{V_n^2} \sum_{k=1}^n \sigma_k^2$$
(3.82)

$$= \int_{|z| > \frac{\varepsilon V_n}{\max_{1 \le k \le n} \sigma_k}} z^2 d\Phi(z)$$
(3.83)

$$\rightarrow 0,$$
 (3.84)

where convergence to zero holds because (3.76) implies divergence of the boundary of integration:  $\frac{\varepsilon V_n}{max_{1\le k\le n}\sigma_k} \to \infty$ . Finally, observing that  $\frac{1}{V_n^2}\sum_{k=1}^n \int_{|x|>\varepsilon V_n} x^2 d|F_k(x) - \Phi_k(x)| \le \frac{1}{V_n^2}\sum_{k=1}^n \int_{|x|>\varepsilon V_n} x^2 dF_k(x) + \frac{1}{V_n^2}\sum_{k=1}^n \int_{|x|>\varepsilon V_n} x^2 d\Phi_k(x)$ , (3.70) follows from Lindeberg's condition and (3.84), so concluding the proof of the implication.

To help an intuitive understanding of the conditions, we note that Lindeberg's condition implies (3.76) and this means that, for large n, the largest variable variance becomes negligible as compared to the total variance. This fact can be interpreted by saying that each variable is infinitesimal if compared to the sum of the others. One the other hand, this infinitesimal behavior is not necessary for the central limit theorem to hold and condition (3.70) is for example satisfied by

$$v_1 \sim G(0,1), v_2 = 0, v_3 = 0, \dots$$
 (3.85)

PROOF OF THEOREM 3.20. For a given *n*, let  $f_k(t)$  and  $\phi_k(t)$  be the characteristic functions of  $v_k/V_n$  and  $z_k/V_n$  (where the  $z_k$ 's are independent and  $G(0, \sigma_k^2)$  distributed),

and let  $f_n(t)$  and  $\phi(t)$  be the characteristic functions of  $\sum_{k=1}^n v_k/V_n$  and  $\sum_{k=1}^n z_k/V_n$  (see Section 2.3 for the notion of characteristic function; note that we have not used the index n in  $\phi(t)$  since  $\sum_{k=1}^n z_k/V_n$  has G(0,1) distribution for any n). Note also that, due to independence,  $f_n(t) = \prod_{k=1}^n f_k(t)$  and  $\phi(t) = \prod_{k=1}^n \phi_k(t)$ .

In the following derivations, the notation  $\int f d(F_k - \Phi_k)$  is short for  $\int f dF_k - \int f d\Phi_k$ . Moreover, we use the following equalities:

$$\int_{\mathbb{R}} dF_k = \int_{\mathbb{R}} d\Phi_k = 1; \quad \int_{\mathbb{R}} x dF_k = \int_{\mathbb{R}} x d\Phi_k = 0; \quad \int_{\mathbb{R}} x^2 dF_k = \int_{\mathbb{R}} x^2 d\Phi_k = \sigma_k^2 \quad (3.86)$$

and the following bounds:

$$\left| e^{\frac{itx}{V_n}} - 1 - it\frac{x}{V_n} + \frac{1}{2}t^2\frac{x^2}{V_n^2} \right| \le t^2\frac{x^2}{V_n^2};$$
(3.87)

$$\left| e^{\frac{itx}{V_n}} - 1 - it\frac{x}{V_n} + \frac{1}{2}t^2\frac{x^2}{V_n^2} \right| \le \frac{1}{6}|t|^3\frac{|x|^3}{V_n^3}$$
(3.88)

(the first bound follows from the Taylor expansion  $e^{\frac{itx}{V_n}} = 1 + it\frac{x}{V_n} - \frac{1}{2}t^2\frac{\xi(x)^2}{V_n^2}$ , with  $|\xi(x)| \le x$ , and the second one from the Taylor expansion  $e^{\frac{itx}{V_n}} = 1 + it\frac{x}{V_n} - \frac{1}{2}t^2\frac{x^2}{V_n^2} - \frac{1}{6}it^3\frac{\xi(x)^3}{V_n^3}$ , again with  $|\xi(x)| \le x$ .) Finally, we also make use of the elementary inequality

$$|\Pi_{k=1}^{n} \alpha_{k} - \Pi_{k=1}^{n} \beta_{k}| \leq \sum_{k=1}^{n} |\alpha_{k} - \beta_{k}|, \qquad (3.89)$$

valid for  $\alpha_k$ 's and  $\beta_k$ 's with  $|\alpha_k|, |\beta_k| \le 1$ . In order to prove (3.89), start with n = 2 and write

$$|\alpha_1\alpha_2 - \beta_1\beta_2| = |\alpha_1\alpha_2 - \beta_1\alpha_2 + \beta_1\alpha_2 - \beta_1\beta_2|$$
(3.90)

$$\leq |(\alpha_1 - \beta_1)\alpha_2| + |\beta_1(\alpha_2 - \beta_2)| \qquad (3.91)$$

$$\leq |\alpha_1 - \beta_1| + |\alpha_2 - \beta_2|. \tag{3.92}$$

The general case is obtained by repeated application of this same inequality. We now have:

$$|f_n(t) - \phi(t)|$$

$$= |\Pi_{t-1}^n f_k(t) - \Pi_{t-1}^n \phi_k(t)|$$
(3.93)
(3.94)

$$\leq \sum_{k=1}^{n} |f_k(t) - \phi_k(t)| \quad (use \ (3.89) \ since \ |f_k(t)|, |\phi_k(t)| \leq 1)$$
(3.95)
(3.95)

$$=\sum_{k=1}^{n}\left|\int_{\mathbb{R}}e^{itx/V_{n}}d(F_{k}-\Phi_{k})\right|$$
(3.96)

$$=\sum_{k=1}^{n} \left| \int_{\mathbb{R}} \left( e^{\frac{itx}{V_n}} - 1 - it \frac{x}{V_n} + \frac{1}{2} t^2 \frac{x^2}{V_n^2} \right) d(F_k - \Phi_k) \right| \quad (use \ (3.86)) \tag{3.97}$$

$$\leq \sum_{k=1}^{n} \int_{|x| \leq \varepsilon V_{n}} \left| e^{\frac{itx}{V_{n}}} - 1 - it \frac{x}{V_{n}} + \frac{1}{2} t^{2} \frac{x^{2}}{V_{n}^{2}} \right| d|F_{k} - \Phi_{k}|$$
(3.98)

$$+\sum_{k=1}^{n}\int_{|x|>\varepsilon V_{n}}\left|e^{\frac{itx}{V_{n}}}-1-it\frac{x}{V_{n}}+\frac{1}{2}t^{2}\frac{x^{2}}{V_{n}^{2}}\right|d|F_{k}-\Phi_{k}|$$
(3.99)

$$\leq \sum_{k=1}^{n} \int_{|x| \leq \varepsilon V_{n}} \frac{1}{6} |t|^{3} \frac{|x|^{3}}{V_{n}^{3}} d|F_{k} - \Phi_{k}| + \sum_{k=1}^{n} \int_{|x| > \varepsilon V_{n}} t^{2} \frac{x^{2}}{V_{n}^{2}} d|F_{k} - \Phi_{k}|$$
(3.100)

(use (3.87) and (3.88)(3.101)

$$\leq \frac{1}{6} |t|^{3} \varepsilon \frac{1}{V_{n}^{2}} \sum_{k=1}^{n} \int_{|x| \leq \varepsilon V_{n}} x^{2} d|F_{k} - \Phi_{k}| + t^{2} \frac{1}{V_{n}^{2}} \sum_{k=1}^{n} \int_{|x| > \varepsilon V_{n}} x^{2} d|F_{k} - \Phi_{k}| (3.102)$$

$$\leq \frac{1}{6}|t|^{3}\varepsilon \frac{1}{V_{n}^{2}} \sum_{k=1}^{n} 2\sigma_{k}^{2} + t^{2} \frac{1}{V_{n}^{2}} \sum_{k=1}^{n} \int_{|x| > \varepsilon V_{n}} x^{2} d|F_{k} - \Phi_{k}|$$
(3.103)

$$= \frac{1}{3} |t|^{3} \varepsilon + t^{2} \frac{1}{V_{n}^{2}} \sum_{k=1}^{n} \int_{|x| > \varepsilon V_{n}} x^{2} d|F_{k} - \Phi_{k}|$$
(3.104)

$$\stackrel{n \to \infty}{\longrightarrow} \frac{1}{3} |t|^3 \varepsilon \quad (use \ (3.70)), \tag{3.105}$$

which, by the arbitrariness of  $\varepsilon$ , implies that

$$|f_n(t) - \phi(t)| \to 0.$$
 (3.106)

The result of the theorem now follows from Theorem 2.11.  $\Box$ 

# 3.5 Weak convergence on $\mathbb{R}$

Weak convergence is a very rich and important topic in measure theory. We discuss only facts related to probability measures on  $\mathbb{R}$  and refer the reader to standard text-

books for more comprehensive treatments.

From Definition 3.1 we know that weak convergence of  $v_n$  to v means that  $\mathbb{E}[f(v_n)] \to \mathbb{E}[f(v)]$  (or, equivalently,  $\int_{\mathbb{R}} f d\mathbb{P}_n \to \int_{\mathbb{R}} f d\mathbb{P}$ , where  $\mathbb{P}_n$  and  $\mathbb{P}$  are the image probabilities of  $v_n$  and v) for any continuous and bounded function  $f : \mathbb{R} \to \mathbb{R}$ . Thus, weak convergence is in fact a property of the image probabilities of the random variables. Making a step towards generality, it should not come as a surprise that weak convergence can also be directly defined for sequences of probability measures.

#### **DEFINITION 3.21** (weak convergence of probability measures on $\mathbb{R}$ ) A

probability measure sequence  $\mathbb{P}_n$  on  $\mathbb{R}$  converges weakly to a probability measure  $\mathbb{P}(\mathbb{P}_n \xrightarrow{w} \mathbb{P})$  if, for any continuous and bounded function  $f : \mathbb{R} \to \mathbb{R}$ , it holds that

$$\int_{\mathbb{R}} f d\mathbb{P}_n \to \int_{\mathbb{R}} f d\mathbb{P}.$$
(3.107)

When  $\mathbb{P}_n \xrightarrow{w} \mathbb{P}$ , we also write  $F_n \xrightarrow{w} F$ , where  $F_n$  and F are the distribution functions associated to  $\mathbb{P}_n$  and  $\mathbb{P}$ .

Thus, weak convergence of  $v_n$  to v can be rephrased by saying that the image probability of  $v_n$  converges weakly to the image probability of v.

Convergence (3.107) is softened by the smoothing properties of integral and, as a consequence, weak convergence can take place among probabilities of very different nature, as the next example illustrates.

**EXAMPLE 3.22** Let  $\mathbb{P}_n$  be the discrete probability with equal mass concentrated in  $\frac{1}{n}, \frac{2}{n}, \ldots, \frac{n}{n}$  (see Figure 3.2.) It follows from Theorem 3.23 below that  $\mathbb{P}_n \xrightarrow{w} \mathbb{P} = \lambda[0, 1]$ , the uniform distribution in [0, 1]. For any n,  $\mathbb{P}_n$  is a discrete distribution and we see that it converges to a distribution of totally different nature:  $\mathbb{P}$  is an absolutely continuous distribution.

The next theorem relates weak convergence on  $\mathbb{R}$  to the behavior of distribution functions.



Figure 3.2: An example of weak convergence:  $F_n$  on the left converges weakly to F, displayed on the right.

**THEOREM 3.23** Let  $\mathbb{P}_n$  and  $\mathbb{P}$  be probability measures on  $\mathbb{R}$ . Then  $\mathbb{P}_n \xrightarrow{w} \mathbb{P}$  if and only if the distribution function  $F_n(x)$  of  $\mathbb{P}_n$  converges to the distribution function F(x) of  $\mathbb{P}$  at every x where F(x) is continuous (in which case we also write  $F_n \xrightarrow{w} F$ ).

In general,  $\mathbb{P}_n \xrightarrow{w} \mathbb{P}$  does not imply that  $F_n(x) \to F(x)$  for those *x* where *F* is not continuous. One example is shown in Figure 3.3 where  $\mathbb{P}_n$  is the concentrated mass in  $\frac{1}{n}$  and  $\mathbb{P}$  is the concentrated mass in 0 (show that  $\mathbb{P}_n \xrightarrow{w} \mathbb{P}$ .) Clearly,  $F_n(0) = 0 \not\to F(0) = 1$ .



Figure 3.3:

#### PROOF.

### $\mathbb{P}_n \xrightarrow{w} \mathbb{P} \Rightarrow F_n(x) \xrightarrow{w} F(x)$ , fat every x where F(x) is continuous)

With x a given point where F is continuous and  $f_+^{\varepsilon}$  and  $f_-^{\varepsilon}$  as represented in Figure 3.4, we have:

$$\limsup_{n \to \infty} F_n(x) \le \limsup_{n \to \infty} \int_{\mathbb{R}} f_+^{\varepsilon} d\mathbb{P}_n = \int_{\mathbb{R}} f_+^{\varepsilon} d\mathbb{P} \le F(x+\varepsilon) \xrightarrow{\varepsilon \downarrow 0} F(x); \quad (3.108)$$

$$\liminf_{n \to \infty} F_n(x) \ge \liminf_{n \to \infty} \int_{\mathbb{R}} f_-^{\varepsilon} d\mathbb{P}_n = \int_{\mathbb{R}} f_-^{\varepsilon} d\mathbb{P} \ge F(x - \varepsilon) \xrightarrow{\varepsilon \downarrow 0} F(x), \quad (3.109)$$



Figure 3.4:

so that  $F_n(x) \to F(x)$ .

 $F_n(x) \to F(x)$  at every x where F(x) is continuous  $\Rightarrow \mathbb{P}_n \xrightarrow{w} \mathbb{P}$ )

Let  $A \in \mathscr{B}(\mathbb{R})$  be a set in  $\mathbb{R}$  with  $\mathbb{P}(\partial A) = 0$  ( $\partial A$  is the boundary of A:  $\partial A = (closure A) \cap (closure A^c)$ ). We want to prove that

$$\mathbb{P}_n(A) \to \mathbb{P}(A). \tag{3.110}$$

Let  $A_0 = interior A$  ("interior A" is the set of points x in A such that  $x \in (a,b)$  for some  $(a,b) \subseteq A$ ). Since  $A_0$  is open, it can be represented as the union of disjoint open intervals:  $A_0 = \bigcup_{k=1}^{\infty} (a_k, b_k)$ . Choose an  $\varepsilon > 0$  and, for each interval  $(a_k, b_k)$  select a subinterval  $(a'_k, b'_k]$  such that  $a'_k$  and  $b'_k$  are points where F(x) is continuous and  $\mathbb{P}(a_k, b_k) \leq P(a'_k, b'_k] + 2^{-k} \cdot \varepsilon$  (since F(x) has at most finitely many discontinuities, such  $a'_k$  and  $b'_k$  certainly exist). Also, fix q such that  $\mathbb{P}(A_0 - \bigcup_{k=1}^q (a_k, b_k)) \leq \varepsilon$ . Now,

$$\mathbb{P}(A_0) = \mathbb{P}(\bigcup_{k=1}^{\infty} (a_k, b_k)) = \sum_{k=1}^{\infty} \mathbb{P}(a_k, b_k) \le \sum_{k=1}^{q} \mathbb{P}(a_k, b_k) + \varepsilon$$
(3.111)

$$\leq \sum_{k=1}^{q} (\mathbb{P}(a'_{k}, b'_{k}] + 2^{-k} \cdot \varepsilon) + \varepsilon = \sum_{k=1}^{q} (F(b'_{k}) - F(a'_{k}) + 2^{-k} \cdot \varepsilon) + (\mathfrak{S}.112)$$

$$\leq \sum_{k=1}^{q} (F(b'_{k}) - F(a'_{k})) + 2\varepsilon = \sum_{k=1}^{q} \lim_{n \to \infty} (F_{n}(b'_{k}) - F_{n}(a'_{k})) + 2\varepsilon \quad (3.113)$$

$$= \lim_{n \to \infty} \sum_{k=1}^{q} \mathbb{P}_n(a'_k, b'_k] + 2\varepsilon \le \liminf_{n \to \infty} \mathbb{P}_n(A_0) + 2\varepsilon, \qquad (3.114)$$

which, due to the arbitrariness of  $\varepsilon$ , gives

$$\mathbb{P}(A_0) \le \liminf_{n \to \infty} \mathbb{P}_n(A_0). \tag{3.115}$$

The same derivation can be applied to  $A_0^c$  (the interior of the complement of A) leading to

$$\mathbb{P}(A_0^c) \le \liminf_{n \to \infty} \mathbb{P}_n(A_0^c). \tag{3.116}$$

Moreover, it holds that

$$\mathbb{P}(\partial A) = 0 \le \liminf_{n \to \infty} \mathbb{P}_n(\partial A), \tag{3.117}$$

where  $\mathbb{P}(\partial A) = 0$  is by assumption and  $0 \leq \liminf_{n \to \infty} \mathbb{P}_n(\partial A)$  is just because a probability cannot be negative. From (3.115), (3.116), and (3.117) it follows that

$$\mathbb{P}(A_0) = \lim_{n \to \infty} \mathbb{P}_n(A_0), \quad \mathbb{P}(A_0^c) = \lim_{n \to \infty} \mathbb{P}_n(A_0^c), \quad 0 = \lim_{n \to \infty} \mathbb{P}_n(\partial A). \tag{3.118}$$

To prove this, suppose that one of these equations is false. Say the first one. Using (3.115) we have

$$\mathbb{P}(A_0) \le \liminf_{n \to \infty} \mathbb{P}_n(A_0) \le \limsup_{n \to \infty} \mathbb{P}_n(A_0), \tag{3.119}$$

and

equality cannot hold throughout since  $\mathbb{P}(A_0) = \liminf_{n \to \infty} \mathbb{P}_n(A_0) = \limsup_{n \to \infty} \mathbb{P}_n(A_0)$  implies that  $\lim_{n\to\infty} \mathbb{P}_n(A_0)$  exists and it is equal to  $\mathbb{P}(A_0)$ , which is the first equation in (3.118). Thus, it must be true that  $\mathbb{P}(A_0) < \limsup_{n \to \infty} \mathbb{P}_n(A_0)$ . Now, this latter fact gives us the possibility of extracting a subsequence  $n_k$  such that  $\mathbb{P}(A_0) < \lim_{k \to \infty} \mathbb{P}_{n_k}(A_0)$ . But this, used together (3.116) and (3.117), leads to an absurd inequality:

$$1 = \mathbb{P}(A_0 \cup A_0^c \cup \partial A) = \mathbb{P}(A_0) + \mathbb{P}(A_0^c) + \mathbb{P}(\partial A)$$
(3.120)

$$< \lim_{k \to \infty} \mathbb{P}_{n_k}(A_0) + \liminf_{n \to \infty} \mathbb{P}_n(A_0^c) + \liminf_{n \to \infty} \mathbb{P}_n(\partial A)$$
(3.121)

$$\leq \lim_{k \to \infty} \mathbb{P}_{n_k}(A_0) + \liminf_{k \to \infty} \mathbb{P}_{n_k}(A_0^c) + \liminf_{k \to \infty} \mathbb{P}_{n_k}(\partial A)$$
(3.122)

$$\leq \liminf_{k \to \infty} (\mathbb{P}_{n_k}(A_0) + \mathbb{P}_{n_k}(A_0^c) + \mathbb{P}_{n_k}(\partial A)) = \liminf_{k \to \infty} \mathbb{P}_{n_k}(A_0 \cup A_0^c \cup \partial A) (3.123)$$
  
= 
$$\liminf_{k \to \infty} 1 = 1.$$
 (3.124)

Thus, our initial assumption that the first equation in (3.118) was false cannot be correct; proceeding similarly for the other equations in (3.118), (3.118) remains proven.

Result (3.110) now immediately follows from (3.118):

$$\lim_{n \to \infty} \mathbb{P}_n(A) \leq \lim_{n \to \infty} \mathbb{P}_n(A_0 \cup \partial A) = \lim_{n \to \infty} \mathbb{P}_n(A_0) + \lim_{n \to \infty} \mathbb{P}_n(\partial A) \quad (3.125)$$

$$= \mathbb{P}(A_0) + 0 \le P(A). \tag{3.126}$$

and

$$\lim_{n \to \infty} \mathbb{P}_n(A) \geq \lim_{n \to \infty} \mathbb{P}_n(A_0) = \mathbb{P}(A_0) = \mathbb{P}_n(A_0) + \mathbb{P}(\partial A) \geq \mathbb{P}(A). \quad (3.127)$$

We now move to prove that  $\mathbb{P}_n \xrightarrow{w} \mathbb{P}$ .

Pick any continuous and bounded  $f : \mathbb{R} \to \mathbb{R}$  and a number M such that |f(x)| < M,  $\forall x$ . Let  $T := \{t \in [-M, M] \text{ such that } P\{f^{-1}(t)\} \neq 0\}$ . T is at most countable.

Consider a partition of [-M, M]

$$-M = t_1 < t_2 < \dots < t_q = M, \tag{3.128}$$

with  $t_k \notin T$ , k = 1, 2, ..., q, and let  $A_k = f^{-1}[t_k, t_{k+1})$  (i.e.  $A_k$  is the set where  $f \in [t_k, t_{k+1})$ ). Since f is continuous,  $f^{-1}(t_k, t_{k+1})$  (the inverse image of the open interval  $(t_k, t_{k+1})$ ) is open, so that  $\partial A_k \subseteq f^{-1}\{t_k\}$  and, being  $t_k \in T$ ,  $\mathbb{P}(\partial A_k) = 0$ , k = 1, ..., q. So,

$$\begin{aligned} \left| \int_{\mathbb{R}} f d\mathbb{P}_{n} - \int_{\mathbb{R}} f d\mathbb{P} \right| & (3.129) \\ &\leq \left| \int_{\mathbb{R}} f d\mathbb{P}_{n} - \sum_{k=1}^{q-1} t_{k} \mathbb{P}_{n}(A_{k}) \right| + \left| \sum_{k=1}^{q-1} t_{k} \mathbb{P}_{n}(A_{k}) - \sum_{k=1}^{q-1} t_{k} \mathbb{P}(A_{k}) \right| + \left| \sum_{k=1}^{q-1} t_{k} \mathbb{P}(A_{k}) - \int_{\mathbb{R}} 3f d\mathbb{P} \right| \\ &\leq 2 \max_{1 \leq k \leq q-1} (t_{k} - t_{k-1}) + \left| \sum_{k=1}^{q-1} t_{k} (\mathbb{P}_{n}(A_{k}) - \mathbb{P}(A_{k})) \right| & (3.131) \end{aligned}$$

$$\stackrel{n \to \infty}{\longrightarrow} \quad 2 \max_{1 \le k \le q-1} (t_k - t_{k-1}), \tag{3.132}$$

where we have used (3.110) when taking the limit. Since the  $t_k$ 's can be selected to be one close to the next at will,  $\max_{1 \le k \le q-1}(t_k - t_{k-1})$  can be made arbitrarily small, leading to the conclusion that  $\int_{\mathbb{R}} f d\mathbb{P}_n \to \int_{\mathbb{R}} f d\mathbb{P}$  as  $n \to \infty$ , and this ends the proof.  $\Box$ 

Suppose we are given a sequence of probability distributions  $F_n$  on  $\mathbb{R}$ . It is true that we can certainly extract a subsequence  $F_{n_k}$  converging weakly to some probability F? The answer is no, as it can be readily verified by taking  $F_n$  to be a step function in n.

The reason why  $F_{n_k}$  fails to exist in the above example is that the mass escapes to infinity. It is an important fact that if such an "escape to infinity" behavior does not take place, then the  $F_n$ 's are packed in such a way that a converging  $F_{n_k}$  certainly exists. This result is stated in the next theorem (the theorem is a particular case of Prokhorov's theorem, which is valid in generic metric spaces).

**THEOREM 3.24 (Helly)** A sequence of probability distribution functions  $F_n$  on  $\mathbb{R}$  is tight if

$$\sup_{n} \int_{|x| \ge M} dF_n(x) \to 0, \quad as \ M \to \infty.$$
(3.133)

Given a tight sequence of probability distribution functions  $F_n$ , there always exists a subsequence  $F_{n_k}$  of  $F_n$  that converges weakly to a limit probability distribution function F.

**PROOF.** Let  $X = \{x_j, j = 1, 2, ...\}$  be a countable dense subset of  $\mathbb{R}$ .

Since  $F_n(x_1) \in [0, 1]$ , we can find a sequence of integers  $n^{(1)} := \{n_1^{(1)}, n_2^{(1)}, \ldots\}$  such that  $F_{n_k^{(1)}}(x_1)$  is convergent. Let us denote by  $\overline{F}(x_1)$  the limiting value. We now restrict attention to sequence  $F_{n_k^{(1)}}$  and extract a subsequence  $n^{(2)} := \{n_1^{(2)}, n_2^{(2)}, \ldots\}$  of  $n^{(1)}$  such that  $F_{n_k^{(2)}}(x_2)$  converges to a limiting value, say  $\overline{F}(x_2)$ . Clearly, being  $F_{n_k^{(2)}}$  a subsequence of  $F_{n_k^{(1)}}$ ,  $F_{n_k^{(2)}}(x_1)$  still converges to  $\overline{F}(x_1)$ . Proceeding along this scheme, we keep constructing nested subsequences of the original sequence  $F_n$  with the property that they converge on an increasing number of points  $x_j$ .

Now, consider the "diagonal" sequence of integers  $n_k := n_k^{(k)}$ . Then, for each  $x_j$  we have

$$F_{n_k}(x_j) \to \bar{F}(x_j). \tag{3.134}$$

The limit probability distribution *F* can now be defined based on  $\overline{F}$ . For all  $x \in \mathbb{R}$  let

$$F(x) = \inf_{x_j > x} \overline{F}(x_j). \tag{3.135}$$

It is not difficult (though it requires some verification) to show that the so constructed F is a probability distribution function (in particular,  $\lim_{x\to\infty} F(x) = 1$  follows from the tightness condition which guarantees that  $F_n(x) > 1 - \varepsilon$ ,  $\forall n$ , for x large enough) and that  $F_{n_k}(x) \to F(x)$  for all x where F is continuous. Then  $F_{n_k} \xrightarrow{w} F$  by virtue of Theorem 3.23 and this completes the proof.

3.5 Weak convergence on  $\ensuremath{\mathbb{R}}$ 

# **Chapter 4**

# **THE PROJECTION THEOREM**

# 4.1 Hilbert spaces

The natural framework to state the projection theorem is a Hilbert space, hence Hilbert spaces are introduced first. A Hilbert space is an inner product space (with a completeness property), and in turn a inner product space is a vector space (with an inner product). We here introduce these notions in a bottom-up fashion, from vector spaces to Hilbert spaces.

**DEFINITION 4.1 (vector space)** A complex vector space (or a vector space over the complex field  $\mathbb{C}$ ) is a triple  $(V, +, \cdot)$  where V is a set and  $+, \cdot$  are two operations. The first operation, called addition, applies two any pair of vectors  $x, y \in V$  and returns a vector  $x + y \in V$ , while the second operation, called scalar multiplication, applies to a pair  $\alpha, x$ , with  $\alpha \in \mathbb{C}$  and  $x \in V$ , and returns a vector  $\alpha \cdot x \in V$ . The following rules apply to the two operations:

(a) x + y = y + x; (b) x + (y+z) = (x+y) + z; (c) there exists a vector 0 (the zero vector) such that x + 0 = x,  $\forall x \in V$ ; (d) to each  $x \in V$ , there corresponds a vector -x (called the opposite vector) such that x + (-x) = 0; (e)  $1 \cdot x = x$ ,  $\forall x \in V$ ; (f)  $\alpha \cdot (\beta \cdot x) = (\alpha\beta) \cdot x$ ; (g)  $\alpha \cdot (x+y) = \alpha \cdot x + \alpha \cdot y$ ; (h)  $(\alpha + \beta) \cdot x = \alpha \cdot x + \beta \cdot x$ .

In the following, we often use the terminology "the vector space V", which is a shorthand used for "the vector space  $(V, +, \cdot)$ " whenever the operations are either clear from the context or their specification is unimportant. We derive some immediate consequences of Definition 4.1.

- (I) Because of (b), we can write x + y + z with no ambiguity;
- (II) there is no other vector y besides the zero vector such that x + y = x,  $\forall x \in V$ . Indeed, take x = 0 in relation x + y = x and write:

$$0 = 0 + y$$
 (4.1)

$$= y + 0 \tag{4.2}$$

$$= y, \qquad (4.3)$$

showing that *y* has to be the zero vector;

(III)  $0 \cdot x = 0, \forall x \in V$  (note that 0 in the left-hand side is the number zero, while 0 in the right-hand side is the vector zero). In fact:

$$y + 0 \cdot x = y + 0 \cdot x + 1 \cdot x + (-x)$$
 (4.4)

$$= y + (0+1) \cdot x + (-x) \tag{4.5}$$

$$= y + 1 \cdot x + (-x) \tag{4.6}$$

$$= y, \quad \forall y \in V, \tag{4.7}$$

so that, by the uniqueness of the zero vector shown in (II),  $0 \cdot x$  has to be 0;

(IV) the opposite of a vector x is unique. Suppose there are two:  $x_1$  and  $x_2$ . Then,

$$x_1 = x_1 + 0 (4.8)$$

$$= x_1 + (x + x_2) \tag{4.9}$$

$$= (x_1 + x) + x_2 \tag{4.10}$$

$$= 0 + x_2$$
 (4.11)

$$= x_2,$$
 (4.12)

so that  $x_1 = x_2$ ; (V)  $-x = -1 \cdot x, \forall x$ . In fact:

$$x + (-1) \cdot x = (1-1) \cdot x \tag{4.13}$$

$$= 0 \cdot x \tag{4.14}$$

$$= 0 (using (III)).$$
 (4.15)

Since the opposite vector is unique (see (IV)), the conclusion follows.

For short, x + (-y) is also written x - y.

**DEFINITION 4.2 (inner product space)** A inner product space is a vector space V where, to each ordered pair of vectors x and y, there is associated a complex number (x, y) called the inner product (or the scalar product) of x and y, with the following properties:

(i)  $(x, y) = \overline{(y, x)}$  overbar denotes complex conjugation); (l) (x+y,z) = (x,z) + (y,z);(m)  $(\alpha \cdot x, y) = \alpha(x,y);$ (n)  $(x,x) \ge 0$ , and (x,x) = 0 implies x = 0.

A vector space over the real field  $\mathbb{R}$  is defined identically to a complex vector space, except that  $\alpha$  and  $\beta$  are real numbers. In this case, the scalar product (x, y) is a real number too. Throughout, we use notations for the complex case, particularization to the real case is straightforward.

**EXAMPLE 4.3** ( $\mathbb{R}^{\mathbf{m}}$ ) For any fixed m, the set  $\mathbb{R}^{m}$  of real valued p-dimensional vectors with addition and scalar multiplication defined in the usual componentwise manner is a real vector space. It becomes a inner product space by the definition:

$$(x,y) = \sum_{k=1}^{m} x_k y_k,$$
(4.16)

where  $x_k [y_k]$  are the components of vector x [y].

**EXAMPLE 4.4** ( $\mathbb{L}^2$ ) Consider the set of square integrable random variables (i.e., random variables  $\xi$  with  $\mathbb{E}[\xi^2] < \infty$ ) defined over a probability space  $(\Omega, \mathscr{F}, \mathbb{P})$ . This set is indicated with  $\mathbb{L}^2$ . With the usual operations of addition and scalar multiplication,  $\mathbb{L}^2$  is a real vector space.

We can try to endow  $\mathbb{L}^2$  with an inner product by defining

$$(\boldsymbol{\xi}, \boldsymbol{\eta}) = \mathbb{E}[\boldsymbol{\xi}\boldsymbol{\eta}]. \tag{4.17}$$

Along this line, however, a difficulty is encountered. Such a difficulty is merely technical, and we prefer to make it explicit to remove any possibility of confusion.

Conditions (i), (l), (m) and the first part of (n)  $((x,x) \ge 0)$  in Definition 4.2 are easy to verify. Instead, a problem turns up with the second part of (n): (x,x) = 0 implies x = 0. In fact, condition  $(\xi, \xi) = \mathbb{E}[\xi^2] = 0$  does not imply that  $\xi = 0$ ,  $\forall \omega \in \Omega$ , the 0 vector in  $\mathbb{L}^2$ ; it only implies that  $\xi = 0$  almost surely. Thus, the condition in (n) of Definition 4.2 that (x,x) = 0 implies x = 0 is violated in this case. Nevertheless, we

insist with the definition  $(\xi, \eta) = \mathbb{E}[\xi \eta]$  with the understanding that  $\mathbb{L}^2$  is not a inner product space in the standard sense of Definition 4.2. It is instead a generalized inner product space where (n) is substituted by:

(n')  $(x,x) \ge 0$ , and (x,x) = 0 for all x in a set Z, where Z is a set containing the 0 vector.

Z has to be interpreted as the set of vectors that are indistinguishable from 0 in the adopted inner product; in  $\mathbb{L}^2$ , it is the set of random variables  $\xi = 0$  almost surely. Adopting this generalized point of view only introduces minor modifications in the theory of inner product spaces and we will make them explicit in relation to the space  $\mathbb{L}^2$  when it will be reconsidered at later stages of this appendix.

It is worth mentioning that a different route can also be adopted to fix this difficulty in an alternative way: instead of viewing  $\mathbb{L}^2$  as a space of random variables, it can be seen as a space of equivalence classes of random variables, where each class contains all variables that differ only on a zero probability set. This corresponds to a coarser-grained viewpoint where one aggregates all variables that are almost surely coincident. This approach, however, introduces some extra complications with measurability issues and we prefer to adopt the first approach where condition (n') substitutes (n). Still, it should be clear that this choice is purely utilitarian, and has no conceptual motivation.

**EXAMPLE 4.5** (C[0,1]) The set of continuous real functions defined on [0,1] with the usual addition and scalar multiplication operations is a real inner product space by the definition

$$(f,g) = \int_0^1 f(r)g(r)dr.$$
 (4.18)

# Geometry of inner product spaces

Quantity

$$\|x\| := \sqrt{(x,x)}$$
(4.19)

is called the norm of vector *x*.

#### **PARALLELOGRAM LAW 4.6**



Figure 4.1: The parallelogram law.

$$||x+y||^2 + ||x-y||^2 = 2||x||^2 + 2||y||^2$$
 (see Figure 4.1) (4.20)

### PROOF. Note first that

$$(x, -y) = (x, -1 \cdot y) (using (V))$$
 (4.21)

$$= -1(y,x)$$
 (using (i) and (m)) (4.22)

$$= -(x,y).$$
 (4.23)

and, similarly,

$$(-x, -y) = (x, y).$$
 (4.24)

By the properties of the inner product and (4.23) and (4.24), we then have

$$\|x+y\|^2 + \|x-y\|^2$$
(4.25)

$$= (x+y,x+y) + (x-y,x-y)$$
(4.26)

$$= ||x||^{2} + ||y||^{2} + (x,y) + (y,x) + ||x||^{2} + ||y||^{2} - (x,y) - (y,x)$$
(4.27)

$$= 2||x||^2 + 2||y||^2.$$
(4.28)

### **SCHWARZ INEQUALITY 4.7**

$$|(x,y)| \le ||x|| ||y||$$
 (see Figure 4.2) (4.29)

PROOF. If ||x|| = 0, then x = 0 (use (n)), so that  $(x, y) = (0, y) = (0 \cdot x, y) = 0$ (x, y) = 0, and (4.29) is true. Suppose  $||x|| \neq 0$ , then





$$0 \leq \left( y - \frac{(y,x)}{\|x\|^2} x, y - \frac{(y,x)}{\|x\|^2} x \right)$$
(4.30)

$$= ||y||^{2} + \frac{|(y,x)|^{2}}{||x||^{2}} - \frac{(y,x)}{||x||^{2}}(x,y) - \frac{(y,x)}{||x||^{2}}(y,x)$$
(4.31)

$$= ||y||^{2} - \frac{|(x,y)|^{2}}{||x||^{2}}, \qquad (4.32)$$

from which (4.29) follows.

# THE TRIANGLE INEQUALITY 4.8



Figure 4.3: The triangle inequality.
PROOF.

$$||x+y||^2 = (x+y,x+y)$$
(4.34)

$$= ||x||^{2} + ||y||^{2} + (x,y) + (y,x)$$
(4.35)

$$\leq ||x||^{2} + ||y||^{2} + 2|(x,y)|$$
(4.36)

$$\leq ||x||^{2} + ||y||^{2} + 2||x|| ||y|| \quad (use \ (4.29)) \tag{4.37}$$

$$= (||x|| + ||y||)^2, \tag{4.38}$$

**DEFINITION 4.9 (orthogonality)** We say that two vectors x and y are orthogonal (and write  $x \perp y$ ) if (x, y) = 0.

#### **PITAGORA'S THEOREM 4.10**

If  $x \perp y$ , then  $||x + y||^2 = ||x||^2 + ||y||^2$  (see Figure 4.4) (4.39)



PROOF.

$$||x+y||^2 = (x+y,x+y)$$
(4.40)

$$= ||x||^{2} + ||y||^{2} + (x, y) + (y, x)$$
(4.41)

$$= ||x||^2 + ||y||^2.$$
(4.42)

**Convergence in inner product spaces** 



$$\boldsymbol{\rho}(\boldsymbol{x}, \boldsymbol{y}) := \|\boldsymbol{x} - \boldsymbol{y}\| \tag{4.43}$$

is called the distance between x and y.

It is easy to verify that  $\rho(\cdot, \cdot)$  in the above definition is indeed a distance, that is, it satisfies the following usual properties of a distance:

$$\rho(x,y) \ge 0$$
 and  $\rho(x,y) = 0$  implies  $x = y$ ;  
 $\rho(x,y) = \rho(y,x)$ ;  
 $\rho(x,y) \le \rho(x,z) + \rho(z,y)$ .

We say that a sequence  $x_n$ ,  $n = 1, 2, \dots$ , converges to x (and write  $x_n \to x$ ) if  $\rho(x_n, x) \to 0$ . It is easy to see that the limit is unique (if  $x_n \to x$  and  $x_n \to y$ , then x = y).

**THEOREM 4.12** If  $x_n \to x$ , then  $(x_n, y) \to (x, y)$ ,  $\forall y$  (this is expressed in words by saying that the inner product is a continuous operator).

PROOF. By assumption  $||x_n - x|| = \rho(x_n, x) \to 0$ . Then, by the Schwarz inequality 4.7,

 $\rightarrow 0.$ 

$$|(x_n, y) - (x, y)| = |(x_n - x, y)|$$
(4.44)

$$\leq \|x_n - x\| \|y\| \tag{4.45}$$

(4.46)

#### **Hilbert spaces**

A sequence  $x_n$  is said to be a Cauchy (or "fundamental") sequence if,  $\forall \varepsilon > 0, \exists N(\varepsilon)$  such that  $\forall p, q \geq N(\varepsilon)$  it holds that  $\rho(x_p, x_q) \leq \varepsilon$ . The following definition is fundamental.

**DEFINITION 4.13 (completeness - Hilbert space)** A inner product space V is complete if every Cauchy sequence converges to a limit point in V. A complete inner product space is called a Hilbert space.  $\Box$ 

In words, a Cauchy sequence is a sequence where any two vectors in its tail are arbitrarily close to each other. If any such sequence converges, we then say that the space is complete.

The following examples illustrate the concept of completeness.

**EXAMPLE 4.14 (Example 4.3 continued -**  $\mathbb{R}^{\mathbf{m}}$  **is complete)** Suppose  $x_n \in \mathbb{R}^m$  is a Cauchy sequence. Hence,  $\forall \varepsilon > 0, \exists N(\varepsilon)$  such that  $\forall p, q \ge N(\varepsilon)$  it holds that  $\rho(x_p, x_q) = ||x_p - x_q|| = \left(\left(\sum_{k=1}^m (x_{pk} - x_{qk})^2\right)^{1/2} \le \varepsilon$  (indexes k identify components). This relation implies  $|x_{pk} - x_{qk}| \le \varepsilon$ , k = 1, 2, ..., m, so that  $x_{nk} \to x_k$  for some  $x_k \in \mathbb{R}$  (this is because the real line is complete, see any text on real analysis), from which  $x_n \to x$ , where x is the vector in  $\mathbb{R}^m$  whose components are  $x_k$ .

**EXAMPLE 4.15 (Example 4.4 continued -**  $\mathbb{L}^2$  **is complete**) In standard inner product spaces where condition (n) in Definition 4.2 holds, if the limit of a sequence exists, then it is unique, and Definition 4.13 requires that such a limit point actually exists for any Cauchy sequence. In the  $\mathbb{L}^2$  space, property (n) has been substituted by property (n') in Example 4.4. Consequently, if  $\xi_n \to \xi$  in  $\mathbb{L}^2$ ,  $\xi_n$  also tends to any other variable  $\xi + \eta$ , with  $\eta = 0$  almost surely, and the uniqueness of the limit is lost. In this context, by completeness it is meant that any Cauchy sequence has (at least) one limit point. If so, all other variables almost surely equal to this limit point will be limit points as well.

Let us then prove that the limit point of a Cauchy sequence always exists. Suppose  $\xi_n$  is a Cauchy sequence in  $\mathbb{L}^2$ , so that  $\forall \varepsilon > 0, \exists N(\varepsilon)$  such that  $\forall p, q \ge N(\varepsilon)$  we have  $E[(\xi_p - \xi_q)^2]^{1/2} \le \varepsilon$ . Then, it is possible to find a sequence of indexes  $n_k$  such that

$$\mathbb{E}[(\xi_{n_{k+1}} - \xi_{n_k})^2]^{1/2} \le \frac{1}{2^k}.$$
(4.47)

To this end, let  $n_1 = N(\frac{1}{2})$ , so that, for any choice of  $n_2$ ,  $\mathbb{E}[(\xi_{n_2} - \xi_{n_1})^2]^{1/2} \leq \frac{1}{2}$  holds. Then, take  $n_2 = N(\frac{1}{4})$  and so on with  $n_3, n_4, \ldots$ . For sequence  $\xi_{n_k}$  we then have

$$\sum_{k=1}^{\infty} \mathbb{E}[|\xi_{n_{k+1}} - \xi_{n_k}|] \leq \sum_{k=1}^{\infty} \mathbb{E}[(\xi_{n_{k+1}} - \xi_{n_k})^2]^{1/2}$$
(4.48)

(use Schwarz inequality (4.7) with x = 1 and  $y = |\xi_{n_{k+1}} - (\xi_{n_k})|$ 

$$\leq \sum_{k=1}^{\infty} \frac{1}{2^k} \quad (use \ (4.47))$$
(4.50)

$$= 1.$$
 (4.51)

Using (4.51), we first prove that  $\xi_{n_k}$  is almost surely convergent to some random variable  $\xi$  and then that  $\xi$  is the  $\mathbb{L}^2$ -limit of the initial sequence  $\xi_n$ , which shows that a Cauchy sequence in  $\mathbb{L}^2$  has limit, so concluding the proof.

Consider the sequence

$$\zeta_k := |\xi_{n_1}| + |\xi_{n_2} - \xi_{n_1}| + \dots + |\xi_{n_k} - \xi_{n_{k-1}}|.$$
(4.52)

This sequence is increasing and, therefore, convergent (either to a finite value or to infinity). The set A where  $\zeta_k \to \infty$  can be expressed as  $A = \bigcap_{p=1}^{\infty} \bigcup_{k=1}^{\infty} \{\zeta_k \ge p\}$ . Since  $\{\zeta_k \ge p\}$  is measurable (i.e., it belongs to  $\mathscr{F}$ ) and a  $\sigma$ -algebra is closed under countable union and intersection, we have that A is measurable too. Moreover, (4.51) implies that  $\mathbb{P}(A) = 0$  (why?). Consider now

$$\xi_{n_k} = \xi_{n_1} + (\xi_{n_2} - \xi_{n_1}) + \dots + (\xi_{n_k} - \xi_{n_{k-1}}).$$
(4.53)

Certainly,  $\xi_{n_k}$  converges everywhere on  $A^c$  (the complement of A) since a sequence is always convergent whenever the corresponding absolute sequence converges to a finite value. Let  $\xi$  be the limit. On A, define  $\xi = 0$ . Then,  $\xi$  is the almost sure limit of  $\xi_{n_k}$ and it is a random variable, i.e., it is measurable (apply Theorem 3.6.)

Finally, we show that  $\xi_n \to \xi$  in  $\mathbb{L}^2$ . Fix  $\varepsilon > 0$  and an integer  $n \ge N(\varepsilon)$  and let

$$\eta_j := inf_{k \ge j} \left(\xi_{n_k} - \xi_n\right)^2.$$
(4.54)

Clearly,  $\eta_j \leq (\xi_{n_j} - \xi_n)^2$ , so that

$$\mathbb{E}[\eta_j] \le \varepsilon^2, \quad \text{for any } j \text{ large enough.}$$
(4.55)

On the other hand,  $\eta_j \uparrow (\xi - \xi_n)^2$  almost surely as  $j \to \infty$  and  $\eta_j \ge 0$ , which, by the monotone convergence Theorem 3.8, implies

$$\mathbb{E}[\eta_j] \uparrow \mathbb{E}\left[ \left(\xi - \xi_n\right)^2 \right]. \tag{4.56}$$

Putting together (4.55) and (4.56) gives  $\mathbb{E}\left[\left(\xi - \xi_n\right)^2\right] \leq \varepsilon^2$ , which, by the arbitrariness of  $\varepsilon$ , implies that  $\xi_n \to \xi$  in  $\mathbb{L}^2$ .

**EXAMPLE 4.16 (Example 4.5 continued -** C[0,1] **is not complete**) *The inner product space* C[0,1] *of Example 4.5 is not complete. In fact the sequence* 

$$f_n(x) = \begin{cases} 1, & x \in \left[0, \frac{1}{2}\right] \\ -(n+1)\left(x - \frac{1}{2}\right) + 1, & x \in \left(\frac{1}{2}, \frac{1}{2} + \frac{1}{n+1}\right] \\ 0, & otherwise, \end{cases}$$
(4.57)

is Cauchy (verify this), but does not converge to a continuous function.

#### **Subspaces**

**DEFINITION 4.17 (subspace)** A subspace S of a vector space V is a subset of V that is itself a vector space, relative to the operations defined in V.  $\Box$ 

It is easy to verify that a subset *S* of a vector space *V* is a subspace if and only if addition and scalar multiplication of vectors in *S* are still in *S*:  $x + y \in S$  if  $x, y \in S$  and  $\alpha \cdot x \in S$  if  $x \in S$ .

**DEFINITION 4.18 (closed subspace)** A subspace S of a Hilbert space is closed if any convergent sequence  $x_n \in S$  has limit in S.

**EXAMPLE 4.19** In  $\mathbb{L}^2$ , the limit point of a convergent sequence is not unique (see Example 4.15). We say that S is closed if any convergent sequence has at least one limit point in S.

The vector space  $\mathbb{L}^2[0,1]$  of measurable functions defined on [0,1] such that  $\int_0^1 f(r)^2 dr < \infty$  is complete (in fact, this is the  $\mathbb{L}^2$  space defined over  $(\Omega, \mathscr{F}, P) = ([0,1], \mathscr{B}[0,1], \lambda)$  and completeness has been proven in Example 4.15). The space C[0,1] of Example 4.16 is a subspace of  $\mathbb{L}^2[0,1]$ , but it is not closed (to verify this, recall that (4.57) does not converge to any function in C[0,1]).

In a Hilbert space H, let  $x^{\perp}$  denote the set of all vectors y orthogonal to x. It is easy to prove that  $x^{\perp}$  is a closed subspace. In fact,  $(y_1 + y_2, x) = (y_1, x) + (y_2, x) = 0$  and  $(\alpha \cdot y, x) = \alpha(y, x) = 0, \forall y_1, y_2, y \in x^{\perp}$ , showing that  $x^{\perp}$  is a subspace. Its closedness is proven by observing that if  $y_n \to y$  and  $y_n \in x^{\perp}$ , then  $0 = \lim_{n \to \infty} (y_n, x) = (y, x)$  (where the last equality follows from Theorem 4.12), that is,  $y \in x^{\perp}$  too.

If *A* is a subset of *H*, by the symbol  $A^{\perp}$  we indicate the set of all vectors orthogonal to every  $x \in A$ . Since  $A^{\perp} = \bigcap_{x \in A} x^{\perp}$ , it is immediate to verify that  $A^{\perp}$  is a closed subspace.

## 4.2 The projection theorem

**THEOREM 4.20 (projection theorem)** Let *S* be a closed subspace of a Hilbert space *H*. Every vector  $x \in H$  has a unique decomposition

$$x = s + z, \tag{4.58}$$

where  $s \in S$  and  $z \in S^{\perp}$  (s is called the projection of x onto S). Moreover, s is the unique vector in S nearest to x:  $||x - s|| = \min_{s' \in S} ||x - s'||$ .

PROOF. To prove the uniqueness of the decomposition, take a hypothetic alternative decomposition  $x = s_1 + z_1$ , with  $s_1 \in S, z_1 \in S^{\perp}$ , and observe that

$$0 = ||x - x||^2$$
(4.59)
$$||z - z - z||^2$$
(4.60)

$$= \|s - s_1 + z - z_1\|^2$$
(4.60)

$$= \|s - s_1\|^2 + \|z - z_1\|^2 \quad (use \ Pitagora's \ Theorem \ 4.10); \qquad (4.61)$$

$$\Rightarrow ||s - s_1||^2 = 0 \text{ and } ||z - z_1||^2 = 0$$
(4.62)

$$\Rightarrow s = s_1 \text{ and } z = z_1, \tag{4.63}$$

so that the two decompositions must coincide.

To prove the existence of the decomposition, note first that if there exists a vector  $s \in S$  at nearest distance from x, then such a s gives the sought decomposition with z := x - s. To show this, we only have to prove that  $x - s \in S^{\perp}$ . Suppose not. Then, take  $y \in S$  such that  $(x - s, y) \neq 0$  and compute

$$\begin{aligned} \left\| x - s - \frac{(x - s, y)}{\|y\|^2} y \right\|^2 &= \left( x - s - \frac{(x - s, y)}{\|y\|^2} y, x - s - \frac{(x - s, y)}{\|y\|^2} y \right) \quad (4.64) \\ &= \|x - s\|^2 - \frac{|(x - s, y)|^2}{\|y\|^2}, \quad (4.65) \end{aligned}$$

which shows that  $s + \frac{(x-s,y)}{\|y\|^2} y \in S$  would be closer to x than s, so contradicting the assumption that s is the vector at nearest distance.

Thus, to prove existence of the decomposition, all we need to show is the existence of a vector  $s \in S$  at nearest distance from x, which is established in the following.

Let  $\delta := \inf_{s' \in S} ||x - s'||$ , and consider a sequence of vectors  $s_n \in S$  such that  $||x - s_n|| \rightarrow \delta$ . We show that  $s_n$  is a Cauchy sequence and that it converges to the sought vector s.

By the parallelogram law 4.6, we have

$$\|s_p - s_q\|^2 = \|(x - s_q) - (x - s_p)\|^2$$
(4.66)

$$= 2\|x - s_q\|^2 + 2\|x - s_p\|^2 - 4\left\|x - \frac{s_q + s_p}{2}\right\|^2.$$
(4.67)

Since  $\frac{s_q+s_p}{2}$  belongs to *S*, the last term  $4 \left\| x - \frac{s_q+s_p}{2} \right\|^2$  is no smaller that  $4\delta^2$ . On the other hand, the first two terms tend to  $2\delta^2$  by construction. Thus,  $\|s_q - s_p\|^2$  is arbitrarily small for any *p* and *q* large enough, that is,  $s_n$  is indeed a Cauchy sequence. Letting *s* be its limit point (which is in *S* by the assumption that *S* is closed), a straightforward application of the triangular inequality now shows that  $\|x - s\| = \delta$ , that is, *s* is at nearest distance:  $\|x - s\| \le \|x - s_n\| + \|s_n - s\| \to \delta + 0 = \delta$ , but  $\|x - s\|$  does not depend on *n* so that  $\|x - s\| \le \delta$ ; on the other hand,  $\|x - s\|$  cannot be smaller than  $\delta$  since  $\delta$  is a lower bound to  $\|x - s'\|$ ,  $\forall s' \in S$ , hence  $\|x - s\| = \delta$ .

Summing up, we have shown that decomposition x = s + z exists and is unique. Moreover, by construction,  $||x - s|| = \delta = \min_{s' \in S} ||x - s'||$ .

Before closing this proof note that the theorem statement contains a very last point: the vector  $s \in S$  nearest to x is unique. This final point is readily established from what we have already proven: suppose for the purpose of contradiction that a second  $s_1 \neq s$  exists in S at nearest distance; then,  $z_1 := x - s_1$  would belong to  $S^{\perp}$ , so providing a second decomposition  $x = s_1 + z_1$ . But this is in contradiction with the already proven uniqueness of the decomposition.

#### Interpretation of the projection theorem

The projection theorem tells us two things.

**1.** If *s* is the projection of *x* onto *S* (i.e. z := x - s is orthogonal to all vectors in *S*), then *s* minimizes the distance of the subspace *S* from *x* (see Figure 4.5). Thus, if we are given the projection, the problem of finding the vector in *S* closest to *x* is automatically solved.

**2.** In addition, the theorem tells us that such a projection actually exists (and is unique) in full generality. The idea in the proof is to construct a sequence  $s_n \in S$  whose distance from *x* tends to the minimal distance  $\delta := \inf_{s' \in S} ||x - s'||$ . Such a sequence tends to accumulate (i.e. it is a Cauchy sequence), so that, by the fundamental closedness assumption of *S*, it converges to a vector *s*, and this *s* is the projection (see Figure 4.6).







Figure 4.6: Construction of *s*.

## **4.3** Applications of the projection theorem

We discuss two applications of the projection theorem. The first one refers to  $\mathbb{R}^m$  and is presented mostly for pedagogical reasons. The second, concerning  $\mathbb{L}^2$ , is of great importance in estimation theory.

#### Application 1: $\mathbb{R}^m$

As we know,  $\mathbb{R}^m$  is a Hilbert space (see Example 4.14). Given  $x \in \mathbb{R}^m$  and  $r \ (r \le m)$  linearly independent vectors  $y_1, y_2, \ldots, y_r$  in  $\mathbb{R}^m$ , consider the problem of finding the vector *s* closest to *x* and belonging to  $S := span\{y_1, y_2, \ldots, y_r\}$ , the subspace linearly spanned by  $y_1, y_2, \ldots, y_r$  (this is certainly a closed subspace, since any finite dimensional subspace is closed).

In matrix notations, vectors  $s \in S$  can be expressed as

$$s = Y\alpha, \tag{4.68}$$

where Y is the matrix  $[y_1 \ y_2 \cdots y_r]$  stacking the  $y_k$  vectors and  $\alpha \in \mathbb{R}^r$  is the vector of

unknowns.

In view of the projection theorem, *s* is the projection of *x* onto *S*. Thus, we have to impose that x - s is orthogonal to all vectors in *S* or, equivalently, to vectors  $y_1, y_2, \ldots, y_r$ :

$$(y_k, x-s) = 0, \quad k = 1, 2, \dots, r.$$
 (4.69)

Using the definition of inner product in  $\mathbb{R}^m$  and relation  $s = Y\alpha$ , we then have

$$y_k^T x = y_k^T Y \alpha, \quad k = 1, 2, \dots, r,$$
 (4.70)

which can be written in a more compact form as

$$Y^T x = Y^T Y \alpha. \tag{4.71}$$

Solving this equation yields  $\alpha$ .

Clearly, the same result can be achieved by direct minimization of ||x - s||. Along this route, we write:

$$\|x - s\|^2 = x^T x + \alpha^T Y^T Y \alpha - 2x^T Y \alpha, \qquad (4.72)$$

whose minimization gives again (4.71).

#### Application 2: $\mathbb{L}^2$

In Example 4.15, we have seen that the space  $\mathbb{L}^2$  of square integrable random variables is complete. Strictly speaking, however, it is not a Hilbert space since condition (n) in Definition 4.2 has been substituted by condition (n') in Example 4.4, so that  $\mathbb{L}^2$  is not a standard inner product space. In particular, this implies that in  $\mathbb{L}^2$  the limit point of a convergent sequence is not unique.

In this context, we say that a subspace *S* is closed if any sequence belonging to *S* that is convergent has at least one limit point in *S*.

By inspecting the proof of the projection Theorem 4.20, we see that the results of this theorem are still valid in the present context with one single amendment: the decomposition is no longer unique. In fact, given the decomposition x = s + z, any  $s_1 \in S$  such that  $s_1 = s$  almost surely gives another valid decomposition  $x = s_1 + (z + s - s_1)$  (note that  $z + s - s_1 \in S^{\perp}$ ). Moreover, no other decompositions are possible besides these. For short, we express this fact by saying that the decomposition is almost surely unique. Similarly, the vector in S nearest to x is not unique and the set of points minimizing the distance is: any  $s_1 \in S$  such that  $s_1 = s$  almost surely.

A notable example of this construction is found when  $S = \mathbb{L}^2(\mathscr{G})$ , the subset of  $\mathbb{L}^2$  formed by all  $\mathscr{G}$ -measurable random variables (here,  $\mathscr{G}$  is any sub  $\sigma$ -algebra of  $\mathscr{F}$ ).

Since the sum of  $\mathscr{G}$ -measurable random variables and their product by a scalar  $\alpha$  is still  $\mathscr{G}$ -measurable,  $\mathbb{L}^2(\mathscr{G})$  is a subspace. It is in fact a closed subspace. Indeed, any Cauchy sequence  $v_n$  in  $\mathbb{L}^2(\mathscr{G})$  is certainly convergent to a point v in  $\mathbb{L}^2$  since  $\mathbb{L}^2$  is complete and, by virtue of Theorem 3.7, we can determine a  $\mathscr{G}$ -measurable limit:  $\bar{v}$  such that  $v_n \to \bar{v}$  in  $\mathbb{L}^2$ .

Thus,  $\mathbb{L}^2(\mathscr{G})$  is a closed subspace and, by the projection theorem, any  $v \in \mathbb{L}^2$  has an almost surely unique projection onto  $\mathbb{L}^2(\mathscr{G})$ . This projection minimizes the distance (i.e. the second order moment) from v among all variables that are  $\mathscr{G}$ -measurable.

For convenience, the results valid for  $\mathbb{L}^2$  are summarized in the following theorem.

**THEOREM 4.21** Let S be a closed subspace of  $\mathbb{L}^2$  (for example,  $S = \mathbb{L}^2(\mathcal{G})$ , the subset of  $\mathbb{L}^2$  formed by all  $\mathcal{G}$ -measurable random variables). Then, given any  $v \in \mathbb{L}^2$ , the projection of v onto S exists, is unique up to almost sure equivalence, that is, given a projection, all other projections are characterized as the set of all random variables in S that are almost surely equal to the given projection. Any projection in the equivalence class minimizes the distance between S and v (i.e.,  $\mathbb{E}[(v - \text{projection of } v)^2] = \min_{s \in S} \mathbb{E}[(v - s)^2]$ ). Moreover, the set of all vectors that minimize the distance coincides with the projection equivalence class, that is no other vector besides those in the projection equivalence class minimizes the distance.  $\Box$ 

## **Chapter 5**

# **CONDITIONAL EXPECTATION AND CONDITIONAL DENSITY**

### 5.1 Conditional expectation

#### **Elementary conditional expectation**

Let *v* be a random variable and let  $\mathscr{D} = \{D_1, D_2, \dots, D_N\}$  be a finite decomposition of the sample space  $\Omega$  (that is,  $\Omega = \bigcup_{k=1}^N D_k$  and  $D_k \cap D_j = \emptyset$  for  $k \neq j$ ) such that  $P(D_k) > 0, k = 1, 2, \dots, N$ . The conditional expectation of *v* with respect to the  $\sigma$ algebra  $\sigma(\mathscr{D})$  generated by  $\mathscr{D}$  (this is the class of all possible unions of sets  $D_k$  plus the empty set) is defined as

$$\mathbb{E}[\nu \mid \boldsymbol{\sigma}(\mathscr{D})] := \sum_{k=1}^{N} \frac{\mathbb{E}[\nu \cdot \mathbf{1}(D_k)]}{P(D_k)} \cdot \mathbf{1}(D_k),$$
(5.1)

where  $1(D_k)$  denotes the indicator function of set  $D_k$ , namely  $1(D_k) = 1$  for  $\omega \in D_k$ and  $1(D_k) = 0$  for  $\omega \notin D_k$ . In Figure 5.1 the conditional expectation of a random variable *v* defined over the sample space  $\Omega = [0, 1]$  is shown. The idea is that the value of  $E[v \mid \sigma(\mathcal{D})]$  over each set  $D_k$  is the mean value of *v* over the same set.

From the above definition, it is clear that

i)  $E[v \mid \sigma(\mathscr{D})]$  is  $\sigma(\mathscr{D})$ -measurable; ii) for any  $A \in \sigma(\mathscr{D})$ ,  $\int_A E[v \mid \sigma(\mathscr{D})]dP = \int_A vdP$ .

The interpretation of i) and ii) is as follows. Because of i), we see that  $E[v \mid \sigma(\mathcal{D})]$  is a simpler random variable than *v* (i.e., it is measurable with respect to a coarser  $\sigma$ -algebra than *v* is). On the other hand, fact ii) says that, from the coarser-grained point of view of  $\sigma(\mathcal{D})$ ,  $\mathbb{E}[v \mid \sigma(\mathcal{D})]$  and *v* are indistinguishable.



Figure 5.1: The conditional expectation in an elementary case.

#### **Definition of conditional expectation**

In probability theory, it is often necessary to take conditional expectation with respect to non-elementary  $\sigma$ -algebras that include zero probability events. We here address such a generalization.

**DEFINITION 5.1 (conditional expectation)** Let v be a random variable defined on a probability space  $(\Omega, \mathscr{F}, \mathbb{P})$  such that  $\mathbb{E}[v]$  exists. Given a  $\sigma$ -algebra  $\mathscr{G} \subseteq \mathscr{F}$ , the conditional expectation of v given  $\mathscr{G}$  is a random variable  $\mathbb{E}[v | \mathscr{G}]$  such that

*i*) 
$$\mathbb{E}[v \mid \mathscr{G}]$$
 *is*  $\mathscr{G}$ -measurable;  
*ii*) for any  $A \in \mathscr{G}$ ,  $\int_A E[v \mid \mathscr{G}] dP = \int_A v dP$ .  $\Box$ 

We shall see below that the conditional expectation exists in full generality. Before that, we prove that the conditional expectation is unique up to almost sure equivalence (i.e., two conditional expectations can only differ from each other on a zero probability set).

Suppose that there are two conditional expectations  $\mathbb{E}[v | \mathscr{G}]_1$  and  $\mathbb{E}[v | \mathscr{G}]_2$  that satisfy i) and ii). Let  $A_+ := \{ \boldsymbol{\omega} : \mathbb{E}[v | \mathscr{G}]_1 > \mathbb{E}[v | \mathscr{G}]_2 \}$ . Then,  $A_+ \in \mathscr{G}$  and  $\int_{A_+} (\mathbb{E}[v | \mathscr{G}]_1 - \mathbb{E}[v | \mathscr{G}]_2) d\mathbb{P} = \int_{A_+} \mathbb{E}[v | \mathscr{G}]_1 d\mathbb{P} - \int_{A_+} \mathbb{E}[v | \mathscr{G}]_2 d\mathbb{P} = \int_{A_+} v d\mathbb{P} - \int_{A_+} v d\mathbb{P} = 0$ , which implies that  $\mathbb{P}(A_+) = 0$  since the integrand  $(\mathbb{E}[v | \mathscr{G}]_1 - \mathbb{E}[v | \mathscr{G}]_2)$  in the first integral is strictly positive over  $A_+$ . Similarly,  $\mathbb{P}(A_-) = \mathbb{P}\{\boldsymbol{\omega} : \mathbb{E}[v | \mathscr{G}]_1 < \mathbb{E}[v | \mathscr{G}]_2\} = 0$  and, hence,  $\mathbb{E}[v | \mathscr{G}]_1 = \mathbb{E}[v | \mathscr{G}]_2$  almost surely.

Each random variable satisfying i) and ii) in Definition 5.1 is a "version" of the conditional expectation. To many purposes, specifying the version is immaterial and it is customary to say: "let us consider the conditional expectation  $\mathbb{E}[v | \mathcal{G}]$  of v given  $\mathcal{G}$ " as a shorthand for "let us consider a version of the conditional expectation  $\mathbb{E}[v | \mathcal{G}]$  of

v given  $\mathscr{G}$ , which we voluntarily do not specify because the specification of the version is unimportant in context under consideration".

The fact that the conditional expectation actually exists is now proven in 3 steps.

#### **STEP 1:** Conditional expectation of a random variable $v \in L^2$ .

Consider the space  $\mathbb{L}^2$  of square integrable random variables, i.e., random variables v with  $E[v^2] < \infty$  ( $\mathbb{L}^2$  is studied in Appendix 4, Examples 4.4 and 4.15). Also, consider  $\mathbb{L}^2(\mathscr{G})$ , the subspace of  $\mathbb{L}^2$  formed by all  $\mathscr{G}$ -measurable random variables. Theorem 4.21 in Appendix 4 proves that the projection of a  $v \in \mathbb{L}^2$  onto  $\mathbb{L}^2(\mathscr{G})$  exists and is unique (up to equivalence). We show that such a projection provides a version of the conditional expectation:

$$\mathbb{E}[v \mid \mathscr{G}] = \text{ projection of } v \text{ onto } \mathbb{L}^2(\mathscr{G}), \tag{5.2}$$

where the right-hand side indicates any projection in the equivalence class. To this end, we need to prove that properties i) and ii) in the Definition 5.1 are fulfilled by the projection.

 $\mathscr{G}$ -measurability of the projection is a direct consequence of the fact that the projection belongs to  $\mathbb{L}^2(\mathscr{G})$ . As for ii), by the definition of projection we have the property that

$$v - ($$
 projection of  $v$  onto  $\mathbb{L}^2(\mathscr{G})) \perp g, \quad \forall g \in \mathbb{L}^2(\mathscr{G}).$  (5.3)

In particular, by taking  $g = 1(A), A \in \mathcal{G}$ , we get:  $\int_A (\text{projection of } v \text{ onto } \mathbb{L}^2(\mathcal{G})) d\mathbb{P} = \int_\Omega (\text{projection of } v \text{ onto } \mathbb{L}^2(\mathcal{G})) \cdot 1(A) d\mathbb{P} = \int_\Omega (v + (\text{projection of } v \text{ onto } \mathbb{L}^2(\mathcal{G}) - v)) \cdot 1(A) d\mathbb{P} = \int_\Omega v \cdot 1(A) d\mathbb{P} = \int_A v d\mathbb{P}$ , that is property ii).

For easy reference, we state the obtained result as a theorem.

**THEOREM 5.2 (conditional expectation of v**  $\in \mathbb{L}^2$ ) *If*  $v \in \mathbb{L}^2$ , *then*  $\mathbb{E}[v | \mathscr{G}]$  *is the projection of v onto the subspace*  $\mathbb{L}^2(\mathscr{G})$  *of all*  $\mathscr{G}$ *-measurable square integrable random variables. Precisely, any projection in the equivalence class is a version of*  $\mathbb{E}[v | \mathscr{G}]$  *and all the versions are obtained by varying the projection in the equivalence class.*  $\Box$ 

#### **STEP 2: Conditional expectation of nonnegative random variables.**

Given  $v \ge 0$ , define the sequence of bounded random variables  $v_n := \min\{v, n\}$ , where n = 1, 2, ... Clearly  $v_n \in \mathbb{L}^2$ , so that  $\mathbb{E}[v_n | \mathscr{G}]$  is defined in Step 1. For any *n*, take a version of  $\mathbb{E}[v_n | \mathscr{G}]$ . The fact that  $v_n$  is nondecreasing implies that  $\mathbb{E}[v_n | \mathscr{G}]$  is almost surely nondecreasing too. Indeed, if *A* is the set where  $\mathbb{E}[v_{n+1} | \mathscr{G}] < \mathbb{E}[v_n | \mathscr{G}]$ , we have:  $0 \ge \int_A (\mathbb{E}[v_{n+1} | \mathscr{G}] - E[v_n | \mathscr{G}]) d\mathbb{P} = \int_A \mathbb{E}[v_{n+1} | \mathscr{G}] d\mathbb{P} - \int_A \mathbb{E}[v_n | \mathscr{G}] d\mathbb{P} =$ 

 $\int_A v_{n+1} d\mathbb{P} - \int_A v_n d\mathbb{P} = \int_A (v_{n+1} - v_n) d\mathbb{P} \ge 0$ , from which we find that equality holds throughout so that  $\mathbb{P}(A) = 0$ . Since  $\mathbb{E}[v_n | \mathscr{G}]$  is nondecreasing almost surely, it converges almost surely (to a finite value or to  $\infty$ .) Where  $\mathbb{E}[v_n | \mathscr{G}]$  does not converge, redefine the limit to be zero. We claim that the limit is a version of the conditional expectation of v given  $\mathscr{G}$ , and we verify this in the following.

**NOTE:** We need to remark the fact that defining  $\mathbb{E}[v | \mathscr{G}] = \lim_{n \to \infty} \mathbb{E}[v_n | \mathscr{G}]$  leaves open the possibility that  $\mathbb{E}[v | \mathscr{G}] = \infty$  on a nonzero probability set even when  $v < \infty$ ,  $\forall \omega \in \Omega$ . To see this, consider the following example: over the probability space  $([0,1], \mathscr{B}[0,1], \lambda)$ , with  $\lambda =$  Lebesgue measure, consider the random variable

$$\begin{cases} 0, & x = 0\\ \frac{1}{x}, & otherwise, \end{cases}$$
(5.4)

and let  $\mathscr{G}$  be the trivial  $\sigma$ -algebra that only contains the empty set  $\emptyset$  and the whole set [0,1]. Being  $\mathscr{G}$  trivial,  $\mathbb{E}[v_n | \mathscr{G}]$  is constant over [0,1] and equal to  $\mathbb{E}[v_n]$ . But  $E[v_n] \to \infty$  as  $n \to \infty$ , showing that  $\mathbb{E}[v | \mathscr{G}] = \infty$ ,  $\forall \omega \in [0,1]$ .

The fact that  $\mathbb{E}[v | \mathscr{G}]$  can possibly be  $\infty$  may seem to pose a difficulty since our definition of random variable (Definition 2.2) and all subsequent developments assume  $v \in \mathbb{R}$ , where  $\mathbb{R}$  does not include  $\pm \infty$ . This difficulty is however easy to circumvent provided that one is ready to work with random variables taking value in  $[-\infty,\infty]$ , the extended set of real numbers. We recall that the arithmetic of  $\mathbb{R}$  is extended to  $[-\infty,\infty]$ with the definitions:  $a + \infty = \infty + a = \infty$  if  $a > -\infty$ , and  $a - \infty = -\infty + a = -\infty$  if  $a < \infty; \infty - \infty$  is not defined.  $a \cdot \pm \infty = \pm \infty \cdot a = \pm \infty$  if a > 0,  $a \cdot \pm \infty = \pm \infty \cdot a = \mp \infty$ if a < 0, and  $a \cdot \pm \infty = \pm \infty \cdot a = 0$  if a = 0. One can easily verify that the commutative and associative laws hold in  $[-\infty,\infty]$  and the distributive law holds in  $[-\infty,\infty]$  as long  $as -\infty$  and  $\infty$  do not appear simultaneously in a sum. The reader is referred to [6] for a more-in-deep treatment of this matter. The definition of random variables and all the subsequent developments can naturally be extended to random variables taking value in  $[-\infty,\infty]$ .

Let us go to verify that properties i) and ii) hold. The measurability property i) follows from Theorem 3.6 applied to sequence  $\mathbb{E}[v_n | \mathscr{G}]$  seen as random variables on  $(\Omega, \mathscr{G}, \mathbb{P})$ (the fact that in Theorem 3.6 convergence takes place to a finite value can be easily generalized to the case at hand here that this value can be  $\infty$ ). As for Property ii), note first that  $\int_A \mathbb{E}[v_n | \mathscr{G}] d\mathbb{P} = \int_A v_n d\mathbb{P}, \forall A \in \mathscr{G}$  (this is Property ii) for random variables in  $\mathbb{L}^2$ ). Then,

$$\int_{A} \mathbb{E}[v \mid \mathscr{G}] d\mathbb{P} = \lim_{n \to \infty} \int_{A} \mathbb{E}[v_n \mid \mathscr{G}] d\mathbb{P}$$
(5.5)

(by the monotone convergence Theorem 3.8) (5.6)

$$= \lim_{n \to \infty} \int_{A} v_n d\mathbb{P}$$
(5.7)

$$= \int_{A} \lim_{n \to \infty} v_n d\mathbb{P}$$
(5.8)

(again by the monotone convergence Theorem 3.8)(5.9)

$$= \int_{A} v d\mathbb{P}.$$
 (5.10)

#### **STEP 3:** Conditional expectation of random variables such that $\mathbb{E}[v]$ exists.

Let  $v^+ := \max\{v, 0\}$  and  $v^- := -\min\{v, 0\}$ . Clearly,  $v = v^+ - v^-$ . We show that a version of the conditional expectation is given by the formula

$$\mathbb{E}[v \mid \mathscr{G}] = \mathbb{E}[v^+ \mid \mathscr{G}] - \mathbb{E}[v^- \mid \mathscr{G}], \qquad (5.11)$$

where in the right-hand side we take any version of the conditional expectation of  $v^+$ and of  $v^-$  and the left-hand side is redefined to be zero where both  $\mathbb{E}[v^+ | \mathscr{G}]$  and  $\mathbb{E}[v^- | \mathscr{G}]$  are  $\infty$ .

We first show that it is not possible that  $\mathbb{E}[v^+ | \mathscr{G}]$  and  $\mathbb{E}[v^- | \mathscr{G}]$  take both value  $\infty$  on a nonzero probability set, so that  $E[v | \mathscr{G}]$  takes on expression (5.11) almost surely: were  $\mathbb{E}[v^+ | \mathscr{G}] = \infty$  on a nonzero probability set, we would then have  $\mathbb{E}[v^+] = \infty$ . Similarly,  $\mathbb{E}[v^- | \mathscr{G}] = \infty$  on a nonzero probability set implies  $\mathbb{E}[v^-] = \infty$ . But this would mean that  $\mathbb{E}[v]$  is not defined (recall that  $\mathbb{E}[v]$  is by definition  $\mathbb{E}[v^+] - \mathbb{E}[v^-]$  provided that not both these expectations are  $\infty$ ), which contradicts our initial assumption that  $\mathbb{E}[v]$  exists. Once we have established that (5.11) holds almost surely, showing that  $\mathbb{E}[v | \mathscr{G}]$  satisfies properties i) and ii) is straightforward (the reader is invited to work out the details).

The following example shows the importance of the assumption that E[v] exists when taking conditional expectation.

**EXAMPLE 5.3** Over  $([0,1], \mathscr{B}[0,1], \lambda)$ , consider the random variable

$$v = \begin{cases} 0, & x = 0 \text{ and } 1\\ \frac{1}{x}, & 0 < x \le 0.5\\ \frac{1}{x-1}, & 0.5 < x < 1, \end{cases}$$
(5.12)

and let  $\mathscr{G} = \{\emptyset, [0,1]\}$ . Here,  $\mathbb{E}[v^+ | \mathscr{G}] = \infty$  in [0,1] and, similarly,  $\mathbb{E}[v^- | \mathscr{G}] = \infty$  in [0,1], so that  $\mathbb{E}[v | \mathscr{G}]$  is not defined. The difficulty arises from the fact that  $\mathbb{E}[v]$  does not exist in this case.

#### **Properties of the conditional expectation**

The following properties are listed without proof. They are all valid almost surely and it is understood that E[v],  $E[v_1]$ , and  $E[v_2]$  are assumed to exist. The reader is referred, among other textbooks, to [7], Chapter 2, for a proof.

1. If  $v_1 \leq v_2$ , then  $\mathbb{E}[v_1 | \mathscr{G}] \leq \mathbb{E}[v_2 | \mathscr{G}]$ ; 2. If  $\alpha$  and  $\beta$  are constants such that  $\alpha \mathbb{E}[v_1] + \beta \mathbb{E}[v_2]$  is defined, then  $\mathbb{E}[\alpha v_1 + \beta v_2 | \mathscr{G}] = \alpha \mathbb{E}[v_1 | \mathscr{G}] + \beta \mathbb{E}[v_2 | \mathscr{G}]$ ; 3. If  $\mathscr{G} = \{\emptyset, \Omega\}$ , then  $\mathbb{E}[v | \mathscr{G}] = \mathbb{E}[v]$ ; 4.  $\mathbb{E}[\mathbb{E}[v | \mathscr{G}]] = \mathbb{E}[v]$ ; 5. If  $\mathscr{G}_1 \subseteq \mathscr{G}_2$ , then  $\mathbb{E}[\mathbb{E}[v | \mathscr{G}_2] | \mathscr{G}_1] = \mathbb{E}[v | \mathscr{G}_1]$ ; 6. Let  $v_1$  be a  $\mathscr{G}$ -measurable random variable and assume that  $\mathbb{E}[v_1v_2]$  exists. Then,  $\mathbb{E}[v_1v_2 | \mathscr{G}] = v_1\mathbb{E}[v_2 | \mathscr{G}]$ .

#### Conditional expectation of v<sub>2</sub> given v<sub>1</sub>

Consider a measurable function  $f : \mathbb{R} \to \mathbb{R}$  and a random variable  $\eta : \Omega \to \mathbb{R}$  and denote by  $\sigma(\eta)$  the  $\sigma$ -algebra generated by  $\eta$  (i.e.,  $\sigma(\eta)$  is the smallest  $\sigma$ -algebra in  $\Omega$  with respect to which  $\eta$  is measurable). Clearly, the random variable  $f(\eta) : \Omega \to \mathbb{R}$ is  $\sigma(\eta)$ -measurable (see Theorem 1.5). It is a remarkable fact that the converse also holds true: if a random variable  $\xi$  is  $\sigma(\eta)$ -measurable, then there exists a measurable function f such that  $\xi = f(\eta)$  for all  $\omega \in \Omega$  (see e.g. [7], Theorem 3, Chapter 2, Section 4, for a proof). For easy reference, we state this fact in the following theorem.

**THEOREM 5.4** Given a random variable  $\eta$ , the set of random variables  $\{f(\eta), with f \text{ measurable function from } \mathbb{R} \text{ to } \mathbb{R}\}$  coincides with the set of  $\sigma(\eta)$ -measurable random variables.

**DEFINITION 5.5** (conditional expectation of  $v_2$  given  $v_1$ ) Given two random variables  $v_1$  and  $v_2$  such that  $\mathbb{E}[v_2]$  exists, consider any version of  $\mathbb{E}[v_2 | \sigma(v_1)]$ . Since this conditional expectation is  $\sigma(v_1)$ -measurable, in view of Theorem 5.4 we can write  $\mathbb{E}[v_2 | \sigma(v_1)] = f(v_1)$ , for some measurable function f, where equality holds  $\forall \omega \in \Omega$ . This function  $f : \mathbb{R} \to \mathbb{R}$  is called a conditional expectation of  $v_2$  given  $v_1$ .

Function *f* in Definition 5.5 is not unique. This can be easily understood by observing that  $f(v_1)$  involves only computing f(x) for values of *x* that correspond to  $v_1(\omega)$  for some  $\omega \in \Omega$ . Thus, changing the value of *f* elsewhere does not affect the value of  $f(v_1)$ . Moreover, if one considers a different version of  $\mathbb{E}[v_2 | \sigma(v_1)]$ , by applying Definition 5.5 one finds that a function *f* such that  $f(v_1) = \mathbb{E}[v_2 | \sigma(v_1)]$  for this new version of the conditional expectation is still a conditional expectation of  $v_2$  given  $v_1$ . This adds an extra degree of freedom in the selection of *f*. It is not difficult to see that, given a conditional expectation *f* of  $v_2$  given  $v_1$ , the set of all conditional expectations is the collection of all measurable functions  $f_1$  that differ from *f* on a set having zero  $\mathbb{P}'_{v_1}$  measure, where  $\mathbb{P}'_{v_1}$  is the image probability induced on  $\mathbb{R}$  by  $v_1$ . Any such function is called a version of the conditional expectation of  $v_2$  given  $v_1$ .

Sometimes, we use the symbol  $\mathbb{E}[v_2 | v_1 = x]$  for f(x). Intuitively,  $\mathbb{E}[v_2 | v_1 = x]$  represents the mean value assumed by  $v_2$  once we know  $v_1$  has the value  $v_1 = x$ .

**EXAMPLE 5.6** Let  $\Omega = \{(0,0), (0,1), (1,0), (1,1)\}$ ,  $\mathscr{F} = all \ subsets \ of \ \Omega$ , and let  $\mathbb{P}$  be specified by  $\mathbb{P}(0,0) = \mathbb{P}(1,1) = \frac{1}{6}$  and  $\mathbb{P}(0,1) = \mathbb{P}(1,0) = \frac{2}{6}$ .

Consider the random variables  $v_1$  and  $v_2$  that assign to the sample outcome (i, j) the value i and j, respectively (see Figure 5.2):

$$v_1(i,j) = i; \quad v_2(i,j) = j.$$
 (5.13)



Figure 5.2:  $v_1$  (•), and  $v_2$  ( $\blacksquare$ ) for Example 5.6.

We have:

$$\mathbb{E}[v_2 \mid \boldsymbol{\sigma}(v_1)] = \begin{cases} \frac{4}{6}, & on \ (0,0) \ and \ (0,1) \\ \frac{2}{6}, & on \ (1,0) \ and \ (1,1), \end{cases}$$
(5.14)

and no other version exists in this case. Moreover, any measurable function  $f : \mathbb{R} \to \mathbb{R}$ with  $f(0) = \frac{4}{6}$  and  $f(1) = \frac{2}{6}$  is a conditional expectation of  $v_2$  given  $v_1$ , as it is readily seen by noting that  $f(v_1) = \mathbb{E}[v_2 \mid \boldsymbol{\sigma}(v_1)]$ .

## 5.2 Conditional density

We here define the conditional density of a random variable  $v_2$  given a second random variable  $v_1$ . Throughout, it is assumed that the two random variables  $v_1$  and  $v_2$  admit joint probability density  $p_{v_1,v_2}(x,y)$ . The reader is referred to standard textbooks for a broader treatment of the concept of conditional measures, among which good references are [1, 2].

**DEFINITION 5.7 (conditional density)** Given a version of the joint probability density  $p_{v_1,v_2}(x,y)$  and a version of the probability density  $p_{v_1}(x)$ , the conditional density of  $v_2$  given  $v_1$  is defined as

$$p_{\nu_{2}|\nu_{1}}(y \mid x) := \begin{cases} \frac{p_{\nu_{1},\nu_{2}}(x,y)}{p_{\nu_{1}}(x)}, & \text{if } p_{\nu_{1}}(x) \neq 0\\ 0, & \text{otherwise.} \end{cases}$$
(5.15)

The conditional density  $p_{v_2|v_1}(y \mid x)$  describes how  $v_2$  distributes under the condition that  $v_1 = x$ .

By varying the version of  $p_{v_1,v_2}$  and  $p_{v_1}$ , different versions of  $p_{v_2|v_1}$  are obtained.

From the definition, it is clear that  $p_{v_2|v_1}$  can be constructed from  $p_{v_1,v_2}$  and  $p_{v_1}$ . On the other hand, if we are given  $p_{v_2|v_1}$ , then  $p_{v_1,v_2}$  and  $p_{v_1}$  cannot be uniquely determined. To see this, note that, if we multiply  $p_{v_1}(x)$  by a function f(x) > 0 such that  $\int p_{v_1}(x)f(x)dx = 1$  and repeat a similar operation for  $p_{v_1,v_2}(x,y)$  so obtaining  $p_{v_1,v_2}(x,y)f(x)$ , then the two densities  $p_{v_1,v_2}(x,y)f(x)$  and  $p_{v_1}(x)f(x)$  are different from the original ones but share with these the same conditional density. As we can see, the indetermination lies in the fact that the division by  $p_{v_1}(x)$  in the definition of the conditional density is a normalization operation that hides the relative probability of different x values.

Since  $p_{v_2|v_1}(y | x)$  describes how  $v_2$  distributes when  $v_1 = x$ , it is an intuitive fact that  $p_{v_2|v_1}(y | x)$  contains a richer knowledge than  $\mathbb{E}[v_2 | v_1 = x]$ , the mean of  $v_2$  for  $v_1 = x$ . The following theorem provides a way to compute  $\mathbb{E}[v_2 | v_1 = x]$  from  $p_{v_2|v_1}(y | x)$ .

**THEOREM 5.8** *Given two random variables*  $v_1$  *and*  $v_2$  *such that*  $\mathbb{E}[v_2]$  *exists, a version of*  $\mathbb{E}[v_2 | v_1 = x]$  *is given by* 

$$\mathbb{E}[v_2 \mid v_1 = x] = \begin{cases} \int_{\mathbb{R}} y p_{v_2 \mid v_1}(y \mid x) dy, & \text{if the integral exists} \\ 0, & \text{otherwise.} \end{cases}$$
(5.16)

**NOTE:** The integral  $\int_{\mathbb{R}} yp_{v_2|v_1}(y \mid x)dy$  in (5.16) is not guaranteed to be defined for all values of x. To understand this, just note that, in correspondence of a given x,  $p_{v_1,v_2}(x,y)$  is substantially arbitrary (the only constraints are due to measurability properties) since a line in  $\mathbb{R}^2$  with fixed coordinate x has zero Lebesgue measure. This arbitrariness can be spent so that the integral in (5.16) is undefined for the selected x.

**PROOF.** Let  $A := \{x : p_{v_1}(x) = 0\}$  and note that, by an application of Theorem 1.12, we have

$$\int_{A \times \mathbb{R}} p_{\nu_1, \nu_2}(x, y) d(x, y) = \int_{\Omega} 1(\nu_1 \in A) dP = \int_A p_{\nu_1}(x) dx = 0,$$
(5.17)

which shows that  $p_{v_1,v_2}(x,y) = 0 \lambda^2 - almost surely$  on  $A \times \mathbb{R}$  (i.e.,  $p_{v_1,v_2}(x,y)$  may be different from zero at most in a zero Lebesgue measure set in  $A \times \mathbb{R}$ ). Consequently,

$$p_{\nu_1,\nu_2}(x,y) = p_{\nu_2|\nu_1}(y \mid x)p_{\nu_1}(x) \quad \lambda^2 - almost \ surely,$$
(5.18)

since the two sides are by definition equal outside  $A \times \mathbb{R}$  while in  $A \times \mathbb{R}$  the right-hand side and the left-hand side are both zero  $\lambda^2 - almost surely$ .

Take now a Borel set *B* in  $\mathbb{R}$ . We want to show that

$$\int_{\mathbb{R}^2} 1(x \in B) y \cdot p_{\nu_2 \mid \nu_1}(y \mid x) p_{\nu_1}(x) d(x, y) = \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} 1(x \in B) y \cdot p_{\nu_2 \mid \nu_1}(y \mid x) p_{\nu_1}(x) dy \right] dx,$$
(5.19)

where it is understood that the inner integral is set to zero at those x's where it is undefined.

To see this, start by noting that the existence of  $\mathbb{E}[v_2] = \mathbb{E}[v_2^+] - \mathbb{E}[v_2^-]$  implies that at least one between  $\mathbb{E}[v_2^+]$  and  $\mathbb{E}[v_2^-]$  is finite. Assuming e.g. that  $\mathbb{E}[v_2^+] < \infty$ , we then have:

$$\infty > \mathbb{E}[1(v_1 \in B) \cdot v_2^+]$$
(5.20)

$$= \int_{\Omega} 1(v_1 \in B) \cdot v_2 \cdot 1(v_2 \ge 0) d\mathbb{P}$$
(5.21)

$$= \int_{\mathbb{R}^2} 1(x \in B) y \cdot 1(y \ge 0) p_{\nu_1, \nu_2}(x, y) d(x, y)$$
(5.22)

$$= \int_{\mathbb{R} \times \{y \ge 0\}} \mathbf{1}(x \in B) y \cdot p_{\nu_2 \mid \nu_1}(y \mid x) p_{\nu_1}(x) d(x, y) \quad (use \ (5.18))$$
(5.23)

$$= \int_{\mathbb{R}} \left[ \int_{\{y \ge 0\}} 1(x \in B) y \cdot p_{v_2|v_1}(y \mid x) p_{v_1}(x) dy \right] dx \quad (use \ Tonelli's \ Theorem \ (5124))$$

$$(5.25)$$

which shows that the inner integral  $\int_{\{y\geq 0\}}$  is less than infinity  $\lambda_x - almost surely$ . (5.19) can now be proven by first rewriting the integral  $\int_{\mathbb{R}^2}$  on the left-hand side as  $\int_{\mathbb{R}\times\{y\geq 0\}} + \int_{\mathbb{R}\times\{y<0\}}$ ; then, by applying Tonelli's theorem to each of these two integrals as in the last step of the previous derivation; and, finally, by noting that the sum of the two integrals can be rewritten as the right-hand-side of (5.19) since the inner integral does not exist (and therefore is redefined to be zero) where both  $\int_{\{y\geq 0\}} = \infty$  and  $\int_{\{y<0\}} = -\infty$ , which only happens on a zero  $\lambda_x$ -measure set.

With the technical results (5.18) and (5.19) in our hands, we can now proceed to write  $\int_{\Omega} 1(v_1 \in B)v_2 dP$  in two different ways and, from a comparison of the results, the theorem conclusion will finally be drawn.

First, we have:

$$\int_{\Omega} 1(v_1 \in B) v_2 d\mathbb{P} \tag{5.26}$$

$$= \int_{\mathbb{R}^2} 1(x \in B) y \cdot p_{\nu_1, \nu_2}(x, y) d(x, y)$$
(5.27)

$$= \int_{\mathbb{R}^2} 1(x \in B) y \cdot p_{\nu_2 \mid \nu_1}(y \mid x) p_{\nu_1}(x) d(x, y) \quad (use \ (5.18)) \tag{5.28}$$

$$= \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} 1(x \in B) y \cdot p_{\nu_2 \mid \nu_1}(y \mid x) p_{\nu_1}(x) dy \right] dx \quad (use \ (5.19)) \tag{5.29}$$

$$= \int_{B} \left[ \int_{\mathbb{R}} y p_{\nu_{2} \mid \nu_{1}}(y \mid x) dy \right] p_{\nu_{1}}(x) dx.$$
 (5.30)

But, we also have:

$$\int_{\Omega} 1(v_1 \in B) v_2 d\mathbb{P}$$
(5.31)

$$= \int_{v_1^{-1}(B)} v_2 d\mathbb{P}$$
(5.32)

$$= \int_{v_1^{-1}(B)} \mathbb{E}[v_2 \mid \boldsymbol{\sigma}(v_1)] d\mathbb{P}$$
(5.33)

$$= \int_{v_1^{-1}(B)} f(v_1) d\mathbb{P} \quad (\text{where } f \text{ is a version of the conditional expectation of } v_2 \text{ giv}(5.34)$$

$$= \int_{\Omega} 1(v_1 \in B) f(v_1) d\mathbb{P}$$
(5.35)

$$= \int_{\mathbb{R}} 1(x \in B) f(x) p_{\nu_1}(x) dx \quad (use \ Theorem \ 1.12)$$
(5.36)

$$= \int_{B} f(x) p_{\nu_{1}}(x) dx.$$
 (5.37)

Comparing (5.30) and (5.37), since *B* is arbitrary we conclude that

$$\int_{\mathbb{R}} y p_{\nu_2 \mid \nu_1}(y \mid x) dy = f(x) \quad P'_{\nu_1} - almost \ surely,$$
(5.38)

that is,  $\int_{\mathbb{R}} y p_{v_2|v_1}(y \mid x) dy$  is a version of  $\mathbb{E}[v_2 \mid v_1 = x]$ .

In this Appendix, we have only considered the conditional expectation in the case of  $\mathbb{R}$ -valued random variables. The extension to multi-dimensional ( $\mathbb{R}^n$ -valued) random variables is straightforward and is defined as the  $\mathbb{R}^n$ -valued random variable whose components are the conditional expectations of the components of the random variable. All other concepts also extend in a natural way to the multi-dimensional case.

## **Chapter 6**

# WIDE-SENSE STATIONARY **PROCESSES**

#### **Definitions and examples** 6.1

Let us consider a sequence of complex-valued random variables  $v_t$  defined over some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , where t runs over positive and negative integers:  $t = \dots, -2, -1, 0, 1, 2, \dots$  The fact that  $v_t$  is "complex-valued" simply means that  $v_t = v_{\mathbb{R},t} + iv_{\mathbb{I},t}$ , where  $v_{\mathbb{R},t}$  and  $v_{\mathbb{I},t}$  are real-valued random variables. The motivation for considering complex - as opposed to real -  $v_t$ 's is that certain derivations become notationally easier in a complex setting.  $v_t$  is called a discrete-time complex-valued stochastic process. When  $v_{I,t} = 0$ , the complex-valued process reduces to a real-valued process.

Throughout this appendix, we assume that variables  $v_t$  are square integrable, i.e.  $\mathbb{E}[|v_t|^2] = \mathbb{E}[v_{\mathbb{R},t}^2 + v_{\mathbb{I},t}^2] < \infty$ . The set of square integrable random variables is indicated with  $\mathbb{L}^2$ . Thus,  $v_t \in \mathbb{L}^2$ , for any *t*.

**DEFINITION 6.1 (wide-sense stationary process)** Process  $v_t$  is said to be wide-sense stationary if its mean is constant:

$$\mathbb{E}[v_t] = \mathbb{E}[v_0], \quad \forall t, \tag{6.1}$$

and the covariance of  $v_t$  and  $v_{t+\ell}$  only depends on the time lag  $\ell$ :

$$\mathbb{E}\left[\left(v_{t+\ell} - E[v_{t+\ell}]\right)\overline{\left(v_t - E[v_t]\right)}\right] = \mathbb{E}\left[\left(v_\ell - E[v_\ell]\right)\overline{\left(v_0 - E[v_0]\right)}\right], \quad \forall t, \ell.$$
(6.2)  
(6.2) Constant denotes complex conjugation.)

(overbar denotes complex conjugation.)

Thus, a wide-sense stationary stochastic process has first and second order statistics that are invariant under time shift.

Function

$$\gamma_{\ell} := \mathbb{E}[(v_{\ell} - E[v_{\ell}])(v_0 - E[v_0])]$$
(6.3)

is called the *auto-covariance function* of process  $v_t$ .

It is not difficult to see that  $\gamma_{\ell}$  is symmetric, i.e.  $\gamma_{\ell} = \overline{\gamma}_{-\ell}$ , and *positive semidefinite*, that is, for any given integer *n* and for any choice of the complex numbers  $a_1, \ldots, a_n$ , it holds that

$$\sum_{k,j=1}^{n} a_k \gamma_{k-j} \overline{a}_j \ge 0.$$
(6.4)

Without any loss of generality, from now on we shall assume that  $\mathbb{E}[v_t] = 0$  (if this is not the case, it is sufficient to subtract the mean from the original stochastic process in order to conform to this assumption). Then,  $\gamma_{\ell} = \mathbb{E}[v_{\ell}\overline{v}_0]$  can be interpreted as the scalar product between the random variables  $v_{\ell}$  and  $v_0$  in the  $\mathbb{L}^2$  space (the reader is referred to Appendix 4 for the notion of scalar product and, particularly, to Example 4.4 for the scalar product in  $\mathbb{L}^2$  - in fact, in Example 4.4 real-valued random variables are considered, but the extension to complex-valued variables presents no difficulties).

The concept of stationary process is now illustrated through examples.

**EXAMPLE 6.2** *Given a real random variable z with*  $\mathbb{E}[z^2] < \infty$ *, let* 

$$v_t = z, \quad \forall t. \tag{6.5}$$

It is immediately seen that  $v_t$  is a wide-sense stationary process. Its realizations are constant functions.  $\Box$ 

**EXAMPLE 6.3 (white process)** A real stochastic process  $v_t$  such that

$$\mathbb{E}[v_t] = 0, \quad \forall t, \tag{6.6}$$

and

$$\mathbb{E}[v_{t+\ell}v_t] = \begin{cases} \sigma^2, & if \ \ell = 0\\ 0, & otherwise, \end{cases}$$
(6.7)

#### is clearly wide-sense stationary.

Such a process is called a white process and its characteristic is that each random variable is uncorrelated with all others. A white process can be equivalently seen as a sequence of orthogonal functions with constant norm in the  $\mathbb{L}^2$  space.

**EXAMPLE 6.4 (one-harmonic process)** Consider the stochastic process  $v_t$  defined through the relation

$$v_t = z e^{-i\omega t} + \bar{z} e^{i\omega t}, \tag{6.8}$$

where  $\omega$  is a fixed frequency belonging to the interval  $[-\pi,\pi]$  and  $z = z_{\mathbb{R}} + iz_{\mathbb{I}}$  is a complex random variable such that:  $\mathbb{E}[z_{\mathbb{R}}] = \mathbb{E}[z_{\mathbb{I}}] = 0$ ,  $\mathbb{E}[z_{\mathbb{R}}^2] = \mathbb{E}[z_{\mathbb{I}}^2] = \sigma^2/4$ ,  $\mathbb{E}[z_{\mathbb{R}}z_{\mathbb{I}}] = 0$ .

In (6.8), process  $v_t$  has been defined through complex quantities for the sake of notational compactness. However, an easy computation shows that process  $v_t$  is in fact real:

$$v_t = (z_{\mathbb{R}} + iz_{\mathbb{I}})(\cos(\omega t) - i\sin(\omega t)) + (z_{\mathbb{R}} - iz_{\mathbb{I}})(\cos(\omega t) + i\sin(\omega t))$$
(6.9)

$$= 2z_{\mathbb{R}}\cos(\omega t) + 2z_{\mathbb{I}}\sin(\omega t)$$
(6.10)

$$= 2\sqrt{z_{\mathbb{R}}^2 + z_{\mathbb{I}}^2 \sin\left(\omega t + atan(z_{\mathbb{R}}/z_{\mathbb{I}})\right)}$$
(6.11)

$$= A\sin(\omega t + \phi), \qquad (6.12)$$

where, in the last equality, we have defined  $A = 2\sqrt{z_{\mathbb{R}}^2 + z_{\mathbb{I}}^2}$  and  $\phi = atan(z_{\mathbb{R}}/z_{\mathbb{I}})$  (here, the appropriate determination for atan has to be taken depending on the sign of  $z_{\mathbb{R}}$  and  $z_{\mathbb{I}}$ ). Expression (6.12) reveals the nature of process  $v_t$ : all its realizations are sinusoids with fixed frequency  $\omega$  and random amplitude A and phase  $\phi$ .

The stationarity of process  $v_t$  can be verified by a direct computation of its mean and auto-covariance function:

$$\mathbb{E}[v_t] = \mathbb{E}\left[ze^{-i\omega t} + \bar{z}e^{i\omega t}\right]$$
(6.13)

$$= \mathbb{E}[z]e^{-i\omega t} + E[\overline{z}]e^{i\omega t}$$
(6.14)

$$= 0, \quad \forall t; \tag{6.15}$$

$$\mathbb{E}[v_{t+\ell}\overline{v_t}] = \mathbb{E}\left[\left(ze^{-i\omega(t+\ell)} + \overline{z}e^{i\omega(t+\ell)}\right)\left(\overline{z}e^{i\omega t} + ze^{-i\omega t}\right)\right]$$
(6.16)

$$= \mathbb{E}[|z|^2]e^{-i\omega\ell} + \mathbb{E}[z^2]e^{-i\omega(2t+\ell)} + \mathbb{E}[\bar{z}^2]e^{i\omega(2t+\ell)} + \mathbb{E}[|z|^2]e^{i\omega\ell}.(6.17)$$

In the latter expression, the expectations are given by

$$\mathbb{E}[|z|^2] = \mathbb{E}[z_{\mathbb{R}}^2 + z_{\mathbb{I}}^2] = \sigma^2/2$$
(6.18)

$$\mathbb{E}[z^2] = \mathbb{E}[z_{\mathbb{R}}^2 - z_{\mathbb{I}}^2 + 2iz_{\mathbb{R}}z_{\mathbb{I}}] = \mathbb{E}[z_{\mathbb{R}}^2] - \mathbb{E}[z_{\mathbb{I}}^2] = 0$$
(6.19)

$$\mathbb{E}[\overline{z}^2] = \mathbb{E}[z_{\mathbb{R}}^2 - z_{\mathbb{I}}^2 - 2iz_{\mathbb{R}}z_{\mathbb{I}}] = \mathbb{E}[z_{\mathbb{R}}^2] - \mathbb{E}[z_{\mathbb{I}}^2] = 0, \qquad (6.20)$$

which, substituted in the expression for  $\mathbb{E}[v_{t+\ell}\overline{v_t}]$ , give

$$\mathbb{E}[v_{t+\ell}\overline{v_t}] = \sigma^2 \cos(\omega \ell), \quad \forall t, \ell.$$
(6.21)

Since this last expression only depends on  $\ell$ , the stationarity of process  $v_t$  follows.  $\Box$ 

**EXAMPLE 6.5 (multi-harmonic process)** The example above can be straightforwardly generalized to the case when the stochastic process is formed by several harmonic components with different frequencies.

Consider

$$v_t = \sum_{k=1}^N \left( z_k e^{-i\omega_k t} + \overline{z}_k e^{i\omega_k t} \right), \tag{6.22}$$

where the  $z_k$ 's satisfy conditions similar to those for z in Example 6.4 and, in addition,  $\mathbb{E}[z_{\mathbb{R},k}z_{\mathbb{R},j}] = \mathbb{E}[z_{\mathbb{R},k}z_{\mathbb{I},j}] = \mathbb{E}[z_{\mathbb{I},k}z_{\mathbb{I},j}] = 0, \forall k \neq j.$ 

Computations entirely similar to those carried out for the case of a single harmonic component show that process  $v_t$  is stationary and that its realizations are formed by the sum of sinusoids with frequencies  $\omega_k$ , k = 1, ..., N. Each sinusoid has random amplitude and phase and the variance of the k-th sinusoidal component is  $\sigma_k^2$ . Moreover, each sinusoidal component is uncorrelated with all others.

In the above example, the stationary stochastic process is the sum of uncorrelated sinusoidal stochastic components. A truly remarkable fact is that this holds true in full generality: any wide-sense stationary stochastic process admits a decomposition in terms of uncorrelated harmonical components. This will be proven later as Theorem 6.10.

## 6.2 Elementary spectral theory of stationary processes

The elementary spectral theory that we present here is not universally applicable, but we prefer to begin with it because it is easy to derive and yet it can be applied to many problems. A general spectral theory is differed to the next Section 6.3.

Assume that  $\gamma_{\ell} \in l^1$  (i.e.  $\sum_{\ell=-\infty}^{\infty} |\gamma_{\ell}| < \infty$ ). This assumption requires that the autocovariance function vanishes fast enough and is not satisfied e.g. in the processes of Examples 6.2, 6.4, and 6.5. Function

$$f(\boldsymbol{\omega}) = \frac{1}{2\pi} \sum_{\ell = -\infty}^{\infty} \gamma_{\ell} e^{-i\boldsymbol{\omega}\ell}$$
(6.23)

(the discrete Fourier transform of  $\gamma_{\ell}$ ) is then pointwise convergent for any  $\omega$  and is called the *spectral density function* of the process.

Clearly,  $f(\omega)$  is periodic of period  $2\pi$ , so that it is enough to regard it as a function defined over  $(-\pi,\pi]$ . Moreover, from the properties of  $\gamma_{\ell}$  it can be seen that  $f(\omega)$  is real and nonnegative.

Given  $\gamma_{\ell}$ , one can compute the spectral density  $f(\omega)$  via (6.23). Viceversa, given  $f(\omega)$ ,  $\gamma_{\ell}$  can be reconstructed by relation

$$\gamma_{\ell} = \int_{(-\pi,\pi]} e^{i\omega\ell} f(\omega) d\omega, \qquad (6.24)$$

(verify this) so that  $\gamma_{\ell}$  and  $f(\omega)$  carry exactly the same information content.

The interpretation of  $f(\omega)$  is that it describes the harmonic content of the stochastic process. A full justification of this interpretation requires a more-in-depth analysis along the lines provided in the next section.

## 6.3 Spectral theory of stationary processes

#### Spectral measure

The spectral measure is a way to prescribe the correlation pattern of a stationary process alternative to  $\gamma_{\ell}$ . Though the spectral measure conveys exactly the same information as  $\gamma_{\ell}$ , it has an extra intuitive appeal because it directly describes the harmonic content of the stationary process. When  $\gamma_{\ell} \in l^1$  as in Section 6.2, the spectral measure has a density and such a density is given by (6.23).

We start with the following fundamental theorem.

**THEOREM 6.6 (Herglotz)** Let  $\gamma_{\ell}$  be a positive semidefinite function (i.e.,  $\gamma_{\ell}$  satisfies (6.4)). Then, there exists a finite measure m on  $((-\pi,\pi], \mathscr{B}(-\pi,\pi])$  such that, for any  $\ell = \ldots, -2, -1, 0, 1, 2, \ldots$ 

$$\gamma_{\ell} = \int_{(-\pi,\pi]} e^{i\omega\ell} dm(\omega).$$
(6.25)

When  $\gamma_{\ell}$  is the auto-covariance function of a wide-sense stationary process  $v_t$ , measure m in Theorem 6.6 is called the *spectral measure* of  $v_t$ . Its distribution function  $F(\omega) = \int_{-\pi}^{\omega} dm(\omega)$  is called the *spectral distribution function*.

PROOF. This proof uses the notion of weak convergence and the reader can find all results used here in Appendix 3.

Define

$$f_n(\omega) := \frac{1}{2\pi n} \sum_{k,j=1}^n e^{-i\omega k} \gamma_{k-j} e^{i\omega j} = \frac{1}{2\pi} \sum_{j=-n+1}^{n-1} \left( 1 - \frac{|j|}{n} \right) \gamma_j e^{-i\omega j}$$
(6.26)

(the second equality follows from a simple computation).

Since  $\gamma_{\ell}$  is positive semidefinite,  $f_n(\omega) \ge 0$ . Next, let

$$F_n(\boldsymbol{\omega}) = \int_{-\pi}^{\boldsymbol{\omega}} f_n(x) dx, \qquad (6.27)$$

for  $\omega \in (-\pi, \pi]$ , while  $F_n(\omega) = 0$  for  $\omega \le -\pi$ , and  $F_n(\omega) = \int_{-\pi}^{\pi} f_n(x) dx$  for  $\omega > \pi$ .  $F_n(\omega)$  is nondecreasing and continuous; moreover, observing that

$$F_n(\pi) = \int_{-\pi}^{\pi} f_n(x) dx = \gamma_0, \quad \forall n,$$
 (6.28)

we conclude that  $F_n/\gamma_0$  is a sequence of probability distribution functions on  $\mathbb{R}$  (see Theorem 2.6 and the comment that follows the statement of this theorem). It is in fact a tight sequence according to the definition of tightness given in Helly's Theorem 3.24 (take  $M = \pi$  in that definition). Thus, due to Theorem 3.24, we conclude that there exists a subsequence  $F_{n_k}/\gamma_0$  which converges weakly to a limit probability distribution function  $F/\gamma_0$ . This distribution function is supported on the closed interval  $[-\pi,\pi]$ . Now, recalling the Definition 3.21 of weak convergence and noting that  $e^{i\omega \ell}$ is a function with continuous and bounded real and complex parts, we obtain

$$\int_{[-\pi,\pi]} e^{i\omega\ell} d\frac{F(\omega)}{\gamma_0} = \lim_{k \to \infty} \int_{[-\pi,\pi]} e^{i\omega\ell} d\frac{F_{n_k}(\omega)}{\gamma_0}$$
(6.29)

$$= \lim_{k \to \infty} \frac{1}{\gamma_0} \int_{[-\pi,\pi]} e^{i\omega \ell} f_{n_k}(\omega) d\omega$$
 (6.30)

$$= \lim_{k \to \infty} \frac{1}{\gamma_0} \int_{[-\pi,\pi]} \frac{1}{2\pi} \sum_{j=-n_k+1}^{n_k-1} \left(1 - \frac{|j|}{n_k}\right) \gamma_j e^{i\omega(-j+\ell)} d\omega 6.31$$

$$\gamma_\ell \qquad (6.22)$$

$$= \frac{n}{\gamma_0}.$$
 (6.32)

Finally, rescale the measure associated to  $F/\gamma_0$  by a factor  $\gamma_0$ . The so-obtained measure is supported on  $[-\pi, \pi]$ , but we can reduce it to a measure supported on  $(-\pi, \pi]$  by transferring the mass in  $-\pi$  to  $\pi$ , and this operation does not change the integral of function  $e^{i\omega\ell}$ . The latter is the measure *m* of the theorem statement, where (6.25) easily follows from (6.32).

A few comments on the spectral measure are in order.

- 1. The spectral measure *m* is defined by means of the auto-covariance function  $\gamma_{\ell}$  only. On the other hand, Herglotz's theorem gives an inversion formula to reconstruct  $\gamma_{\ell}$  from *m*. This shows that *m* and  $\gamma_{\ell}$  carry the same content of information: they both completely define the correlation pattern of the stationary process;
- 2. the spectral measure is unique. This claim requires a simple proof.

Suppose there are two spectral measures with distribution functions  $F_1$  and  $F_2$  (we define  $F_1(\omega) = F_2(\omega) = 0$  for  $\omega \le -\pi$  and  $F_1(\omega) = F_2(\omega) = \gamma_0$  for  $\omega > \pi$ ). Then,

$$\int_{(-\pi,\pi]} e^{i\omega\ell} dF_1(\omega) = \gamma_\ell = \int_{(-\pi,\pi]} e^{i\omega\ell} dF_2(\omega).$$
(6.33)

Given an arbitrary  $\overline{\omega} \in (-\pi, \pi]$ , consider the sequence of functions

$$g_n(\boldsymbol{\omega}) = \begin{cases} 1, & \text{in } (-\pi, \overline{\boldsymbol{\omega}}] \\ 1 - n(\boldsymbol{\omega} - \overline{\boldsymbol{\omega}}), & \text{in } (\overline{\boldsymbol{\omega}}, \overline{\boldsymbol{\omega}} + \frac{1}{n}] \\ 0, & \text{in } (\overline{\boldsymbol{\omega}} + \frac{1}{n}, \pi]. \end{cases}$$
(6.34)

Since every bounded continuous function can be uniformly approximated on  $(-\pi, \pi]$  by trigonometric polynomials (see e.g. [6]), from (6.33) we have

$$\int_{(-\pi,\pi]} g_n(\omega) dF_1(\omega) = \int_{(-\pi,\pi]} g_n(\omega) dF_2(\omega).$$
(6.35)

Sending  $n \to \infty$ , this last equation gives  $F_1(\overline{\omega}) = F_2(\overline{\omega})$ , which, owing to the arbitrariness of  $\overline{\omega}$ , yields the desired result.

Notice also that the uniqueness of the spectral measure implies that one cannot find two subsequences,  $F'_{n_k}/\gamma_0$  and  $F''_{n_k}/\gamma_0$ , of  $F_n/\gamma_0$  that are weakly convergent to two distinct limit distribution functions.

#### Spectral density

Suppose there exists a measurable function f defined on  $(-\pi, \pi]$  whose integral returns the spectral distribution function  $F: F(\omega) = \int_{-\pi}^{\omega} f(x) dx$ . Then, f is called the *spectral density function* of the process. In this case, F is  $\lambda$ -*almost surely* differentiable and fis  $\lambda$ -*almost surely* the derivative of F (this is the "fundamental theorem of calculus", see e.g. Theorem 7.20 in [6]).

When  $\gamma_{\ell} \in l^1$ , f is given by (6.23) (and this justifies our calling "spectral density function" the function in (6.23)), a fact that is proven in the next theorem.

**THEOREM 6.7** *If*  $\gamma_{\ell} \in l^1$ *, then* 

$$f(\boldsymbol{\omega}) = \frac{1}{2\pi} \sum_{\ell = -\infty}^{\infty} \gamma_{\ell} e^{-i\boldsymbol{\omega}\ell}$$
(6.36)

is the spectral density function of the process.

PROOF. Recall that

$$F_n(\omega) = \int_{-\pi}^{\omega} \frac{1}{2\pi} \sum_{j=-n+1}^{n-1} \left(1 - \frac{|j|}{n}\right) \gamma_j e^{-ixj} dx.$$
(6.37)

Since  $\gamma_{\ell} \in l^1$ , the integrand can be uniformly (with respect to *n*) bounded as follows:

$$\left|\frac{1}{2\pi}\sum_{j=-n+1}^{n-1}\left(1-\frac{|j|}{n}\right)\gamma_{j}e^{-ixj}\right| \leq \frac{1}{2\pi}\sum_{j=-\infty}^{\infty}|\gamma_{j}| < \infty.$$
(6.38)

Thus, by the dominated convergence Theorem 3.9 (in fact, here we are integrating with respect to a finite measure instead of a probability measure as in Theorem 3.9, but this difference can be leveled by a rescaling factor), we obtain

$$\lim_{n \to \infty} \int_{-\pi}^{\omega} \frac{1}{2\pi} \sum_{j=-n+1}^{n-1} \left( 1 - \frac{|j|}{n} \right) \gamma_j e^{-ixj} dx = \int_{-\pi}^{\omega} \lim_{n \to \infty} \frac{1}{2\pi} \sum_{j=-n+1}^{n-1} \left( 1 - \frac{|j|}{n} \right) \gamma_j e^{-ixj} dx$$

$$= \int_{-\pi}^{\omega} \frac{1}{2\pi} \sum_{j=-\infty}^{\infty} \gamma_j e^{-ixj} dx, \qquad (6.40)$$

showing that the sequence  $F_n(\omega)$  converges for any  $\omega$  to  $F(\omega) = \int_{-\pi}^{\omega} \frac{1}{2\pi} \sum_{j=-\infty}^{\infty} \gamma_j e^{-ixj} dx$ . Thus,  $F(\omega)/\gamma_0$  is the weak limit of  $F_n(\omega)/\gamma_0$  and so it is the spectral distribution function and  $\frac{1}{2\pi\gamma_0} \sum_{j=-\infty}^{\infty} \gamma_j e^{-i\omega j}$  is its density function. This concludes the proof.

The above proof relies on the fact that a function  $\gamma_{\ell} \in l^1$  has such a thin tail that the limit for  $n \to \infty$  of the integral  $\int_{-\pi}^{\pi} \frac{1}{2\pi} \sum_{j=-n+1}^{n-1} \left(1 - \frac{|j|}{n}\right) \gamma_j e^{-ixj} dx$  can be computed by sending to infinity the summation under the sign of integration. While we have presented this proof of Theorem 6.7 because it is instructive, we also note that an indirect, and quicker, proof can be obtained along the same lines as in point 2 after Theorem 6.6: Theorem 6.6 gives

$$\gamma_{\ell} = \int_{(-\pi,\pi]} e^{i\omega\ell} dF(\omega), \qquad (6.41)$$

where F is the probability distribution function of m, while formula (6.24) gives

$$\gamma_{\ell} = \int_{(-\pi,\pi]} e^{i\omega\ell} f(\omega) d\omega, \qquad (6.42)$$

from which

$$\int_{(-\pi,\pi]} e^{i\omega\ell} dF(\omega) = \int_{(-\pi,\pi]} e^{i\omega\ell} f(\omega) d\omega.$$
(6.43)

Hence, following the same rationale as in point 2 after Theorem 6.6, one can use functions  $g_n(\omega)$  in (6.34) to conclude that

$$F(\bar{\boldsymbol{\omega}}) = \int_{-\pi}^{\bar{\boldsymbol{\omega}}} f(x) dx.$$
(6.44)

**EXAMPLE 6.8 (Example 6.3 continued)** For the white process of Example 6.3,  $F_n(\omega)$  given in (6.37) can be computed as follows

$$F_{n}(\omega) = \int_{-\pi}^{\omega} \frac{1}{2\pi} \sum_{j=-n+1}^{n-1} \left(1 - \frac{|j|}{n}\right) \gamma_{j} e^{-ixj} dx = \int_{-\pi}^{\omega} \frac{1}{2\pi} \sigma^{2} dx \qquad (6.45)$$

$$= \frac{1}{2\pi}\sigma^2(\omega+\pi) = F(\omega), \quad for \ \omega \in (-\pi,\pi].$$
(6.46)

Taking derivative with resepct to  $\omega$ , we find the spectral density to be  $f(\omega) = \frac{1}{2\pi}\sigma^2$ ,  $\lambda$ -almost surely. The same result can be obtained by relation  $f(\omega) = \frac{1}{2\pi}\sum_{\ell=-\infty}^{\infty}\gamma_{\ell}e^{-i\omega\ell} = \frac{1}{2\pi}\sigma^2$ .

**EXAMPLE 6.9 (Example 6.2 continued)** For the process of Example 6.2, assume  $\mathbb{E}[z] = 0$  and  $\mathbb{E}[z^2] = 1$ . Then, the auto-covariance function is  $\gamma_{\ell} = 1$ ,  $\forall \ell$ , and is not in  $l^1$ . In this case, the series in (6.23) is not convergent (consider for example  $\omega = 0$ ). Still, convergence holds in a weak sense as indicated in the proof of Herglotz Theorem 6.6.



Figure 6.1:

Figure 6.1 displays the functions  $F_n(\omega)$  given by (6.37) for some values of n. It can be noted that  $F_n(\omega)$  seems to converge to the step function with a step in 0 of hight 1. This is in fact true and  $F_n(\omega)$  converges to this step function for any  $\omega \neq 0$ , while it holds that  $F_n(0) = 1/2$ ,  $\forall n$  (verifying this claim requires some lengthy computations and the reader can go through the calculations along the following line: first show that  $F_n(\omega) = \frac{1}{2\pi} \sum_{j=-n+1}^{-1} \left(1 - \frac{|j|}{n}\right) \frac{i}{j} e^{-i\omega j} + \frac{1}{2\pi} \sum_{j=1}^{n-1} \left(1 - \frac{|j|}{n}\right) \frac{i}{j} e^{-i\omega j} + \frac{\omega}{2\pi} + \frac{1}{2}$ ; then, notice that term  $\frac{\omega}{2\pi}$  can be seen as the restriction to  $(-\pi,\pi]$  of a periodic saw-tooth function with period  $2\pi$  and compute the Fourier expansion of this function; finally, after substituting the Fourier expansion for  $\frac{\omega}{2\pi}$  in the expression of  $F_n(\omega)$ , one can recognize that the so-obtained expansion for  $F_n(\omega)$  tends to the expansion for the step function).

What is F, the spectral distribution function, in this case? It is the step function with a step in 0 of hight equal to 1. This distribution function is not absolutely continuous and the spectral density function does not exist in this case. The spectral measure has concentrated mass in 0.

We conclude with a technical remark relative to a point that may have attracted the reader's attention. In this example, in order for F to be a distribution function, F(0) must be equal to 1, for, otherwise, F would not be continuous on the right in  $\omega = 0$ . If we go back to the proof of Helly's Theorem 3.24 (Helly's theorem is used in the proof of Herglotz Theorem 6.6), we see that F(x) is defined for any x as  $\inf_{x_j > x} \overline{F}(x_j)$ . By this definition, one indeed has that F(0) = 1.

#### A remark on the usefulness of the spectral distribution

When  $\gamma_{\ell} \in l^1$ , the tail of  $\gamma_{\ell}$  is vanishing fast enough and  $\frac{1}{2\pi} \sum_{\ell=-n}^{n} \gamma_{\ell} e^{-i\omega\ell}$  gently converges as  $n \to \infty$ , allowing us to work with its limit  $\frac{1}{2\pi} \sum_{\ell=-\infty}^{\infty} \gamma_{\ell} e^{-i\omega\ell}$ . This situation, however, does not cover all possible cases and the  $\gamma_{\ell}$  can as well have a powerful tail so that taming the limit behavior of  $\frac{1}{\sqrt{2\pi}} \sum_{\ell=-n}^{n} \gamma_{\ell} e^{-i\omega\ell}$  becomes difficult. The idea behind the construction of the spectral distribution function is to control the "roughness" of  $\frac{1}{2\pi} \sum_{\ell=-n}^{n} \gamma_{\ell} e^{-i\omega\ell}$  by the smoothing properties of integration. The integrated function converges to the spectral distribution function in a weak sense only, but this convergence is strong enough to secure the fundamental inversion formula (6.25).

#### Spectral representation of stationary processes

We now introduce an alternative representation of a wide-sense stationary process where the process is viewed as a stochastic integral with respect to an orthogonal stochastic measure. This representation shows that any stationary process can be interpreted as the sum of elementary harmonic components, so generalizing the situation in Example 6.5. This deep-seated result sheds new light on the very nature of wide-sense stationary processes. Moreover, it lends a new interpretation of the spectral measure mas a quantifier of the harmonic content of the stationary process.

Let us consider the space  $\mathbb{L}^2(m)$  of the complex-valued, measurable and square integrable functions defined on  $((-\pi, \pi], \mathscr{B}(-\pi, \pi], m)$ , where *m* is the spectral measure associated to process  $v_t$ . This space is entirely similar to the space  $\mathbb{L}^2$  studied in Examples 4.4 and 4.15 in Appendix 4, to which we refer the reader for definitions and explanation (in Appendix 4, real-valued functions are considered. Extending the results therein to the present complex-valued setting with a measure *m* presents no difficulties and the reader is invited to work out the details). Here, we merely recall that  $\mathbb{L}^2(m)$  is a vector space that can be endowed with a generalized inner product by the definition

$$(f,g) = \int_{(-\pi,\pi]} f(\boldsymbol{\omega})\overline{g}(\boldsymbol{\omega})dm(\boldsymbol{\omega}).$$
(6.47)

This inner product is "generalized" since (f, f) = 0 does not imply that f = 0, it simply yields f = 0 *m*-almost surely. Importantly,  $\mathbb{L}^2(m)$  is complete, see Example 4.15.

Let  $\mathbb{L}^2_0(m) \subseteq \mathbb{L}^2(m)$  be the linear subspace spanned by  $e^{i\omega k}$ , k = ..., -2, -1, 0, 1, 2, ...(i.e., each element in  $\mathbb{L}^2_0(m)$  is simply a linear combination of functions  $e^{i\omega k}$  of the form  $\sum_{k \in K} \alpha_k e^{i\omega k}$ , where *K* is any finite set of integers and  $\alpha_k$  are complex numbers). The closure of  $\mathbb{L}^2_0(m)$  coincides with  $\mathbb{L}^2(m)$  itself (see e.g. [6]).

Next, introduce in  $\mathbb{L}^2$ , the space of square integrable random variables, the linear subspace  $\mathbb{L}^2_o(v)$  of the variables spanned by  $v_k$  (i.e.,  $\sum_{k \in K} \alpha_k v_k$  with *K* finite and  $\alpha_k$  complex) and denote by  $\mathbb{L}^2(v)$  its closure. In general,  $\mathbb{L}^2(v) \neq \mathbb{L}^2$ .

We want to establish a one-to-one correspondence *T* between  $\mathbb{L}^2(m)$  and  $\mathbb{L}^2(v)$ . To be precise, *T* is one-to-one up to equivalence, i.e. we identify functions of  $\mathbb{L}^2(m)$  which are *m*-almost surely equal and random variables of  $\mathbb{L}^2(v)$  which are *P*-almost surely equal.

To this end, define first the following correspondence between elements of  $\mathbb{L}^2_0(m)$  and elements of  $\mathbb{L}^2_0(v)$ :

$$\sum_{k \in K} \alpha_k e^{i\omega k} \leftrightarrow \sum_{k \in K} \alpha_k v_k.$$
(6.48)

This correspondence is an isometry, that is, it preserves the inner product. In fact:

$$\left(\sum_{k\in K}\alpha_{k}e^{i\omega k},\sum_{j\in J}\beta_{j}e^{i\omega j}\right) = \int_{(-\pi,\pi]}\left(\sum_{k\in K}\alpha_{k}e^{i\omega k}\right)\left(\sum_{j\in J}\overline{\beta}_{j}e^{-i\omega j}\right)dm(\omega)(6.49)$$
$$= \sum\sum_{i\in J}\sum_{j\in J}\alpha_{i}\overline{\beta}_{j}\int_{-\infty}e^{i\omega(k-j)}dm(\omega)(6.50)$$

$$\sum_{k \in K} \sum_{j \in J} \alpha_k \overline{\beta}_j \int_{(-\pi,\pi]} e^{i\omega(k-j)} dm(\omega)$$
(6.50)

$$= \sum_{k \in K} \sum_{j \in J} \alpha_k \overline{\beta}_j \gamma_{k-j} \quad (by \ Theorem \ 6.6) \tag{6.51}$$

$$= \sum_{k \in K} \sum_{j \in J} \alpha_k \overline{\beta}_j E[v_k \overline{v}_j]$$
(6.52)

$$= \left(\sum_{k\in K} \alpha_k v_k, \sum_{j\in J} \beta_j v_j\right).$$
(6.53)

Next, consider an  $f \in \mathbb{L}^2(m)$ . Since  $\mathbb{L}^2(m)$  is the closure of  $\mathbb{L}^2_0(m)$ , there is a sequence  $f_n \in \mathbb{L}^2_0(m)$  that converges to f. Let  $z_n$  be the sequence of random variables in  $\mathbb{L}^2_0(v)$ 

corresponding to  $f_n$ . Since the correspondence is an isometry, it is immediate to verify that  $z_n$  is a Cauchy sequence, so that, in view of the completeness of  $\mathbb{L}^2$ , it converges to some limit point in  $\mathbb{L}^2(v)$ . There is an evident converse to this construction: if  $z \in \mathbb{L}^2(v)$ , then one can find a sequence in  $\mathbb{L}^2_0(v)$  converging to z and the corresponding sequence in  $\mathbb{L}^2_0(m)$  converges to a function  $f \in \mathbb{L}^2(m)$ .

By definition, we let  $T : f \leftrightarrow z$ . It is easy to verify that this correspondence between  $\mathbb{L}^2(m)$  and  $\mathbb{L}^2(v)$  is one-to-one, linear and isometric. The construction is illustrated in Figure 6.2.



Figure 6.2: The correspondence *T*.

We next introduce the notion of stochastic integral. For any Borel set  $B \in \mathscr{B}(-\pi,\pi]$ , denote by T(B) the random variable corresponding to 1(B) (the indicator function of set *B* equal to 1 on *B* and 0 elsewhere) in the isometry *T*. Note that  $\mathbb{E}[T(B_1)\overline{T}(B_2)] = 0$  whenever  $B_1 \cap B_2 = \emptyset$ . This is immediately verified by noting that  $\mathbb{E}[T(B_1)\overline{T}(B_2)] = \int_{(-\pi,\pi]} 1(B_1)1(B_2)dm(\omega) = \int_{(-\pi,\pi]} 0 dm(\omega) = 0$ . Function  $T : \mathscr{B}(-\pi,\pi] \to \mathbb{L}^2(v)$  is called an orthogonal stochastic measure.

We first define the stochastic integral for elementary functions. Given a finite set of N disjoint Borel sets  $B_k$ , k = 1, ..., N, and N complex numbers  $\alpha_k$ , k = 1, ..., N, the integral of the elementary function  $f = \sum_{k=1}^{N} \alpha_k 1(B_k)$  is simply defined as  $\sum_{k=1}^{N} \alpha_k \cdot T(B_k)$ . Here, one should note that the stochastic integral of f is nothing but the random variable that corresponds to f in the isometry T.

Next, the definition is extended to any function  $f \in L^2(m)$  as follows. Take a sequence  $f_n$  of elementary functions that converges to f in  $L^2(m)$ . Since the integrals of the  $f_n$ 's are the random variables corresponding to  $f_n$  in the isometry, such integrals form

a Chauchy sequence and therefore converge to a limit point in  $\mathbb{L}^2(v)$ . This limit is by definition the integral of f and it is denoted by  $\int_{(-\pi,\pi]} f(\omega) dT(\omega)$ . Again, the integral of f is nothing but the random variable corresponding to f in the isometry T.

The notion of stochastic integral can now be applied to the functions  $e^{i\omega t}$ . By construction, we obtain  $\int_{(-\pi,\pi]} e^{i\omega t} dT(\omega) = v_t$ .

We have proved the following result.

**THEOREM 6.10** There exists an orthogonal stochastic measure T on  $((-\pi,\pi], \mathscr{B}(-\pi,\pi])$  such that

$$v_t = \int_{(-\pi,\pi]} e^{i\omega t} dT(\omega).$$
(6.54)

Moreover,  $\mathbb{E}[|T(B)|^2] = m(B), \forall B \in \mathscr{B}(-\pi, \pi].$ 

#### A new interpretation of stationary processes

The above definition of stochastic integral delivers a new interpretation of a stationary process that points directly to its inborn structure and provides an insightful standpoint in many applications.

Consider equation  $v_t = \int_{(-\pi,\pi]} e^{i\omega t} dT(\omega)$ . Let us partition the interval  $(-\pi,\pi]$  in a large, though finite, number of subintervals of equal length:  $(-\pi,\pi] = (-\pi,-\pi + 2\pi\frac{1}{N}] \cup (-\pi + 2\pi\frac{1}{N}, -\pi + 2\pi\frac{2}{N}] \cup \cdots \cup (\pi - 2\pi\frac{1}{N}, 2\pi] = \bigcup_{k=1}^{N} B_k$ . Then, functions  $e^{i\omega t}$ ,  $t \leq T$ , can be approximated by  $\sum_{k=1}^{N} e^{i\omega_k t} 1(B_k)$  (when *t* becomes too large, beyond *T*, sum  $\sum_{k=1}^{N} e^{i\omega_k t} 1(B_k)$  fails to approximate  $e^{i\omega t}$  because  $e^{i\omega t}$  oscillates too fast as a function of  $\omega$ ). When this happens, one needs to introduce a more fine-grained partition of  $(-\pi,\pi]$ . Correspondingly,  $v_t, t \leq T$ , can be approximated by the stochastic integral of this latter function, leading to

$$v_t \approx \sum_{k=1}^N e^{i\omega_k t} T(B_k), \quad t \le T.$$
(6.55)

This expression delivers an interpretation of process  $v_t$  which is very useful for an intuitive understanding of its structure:

A stationary process  $v_t$  is given by the linear combination of uncorrelated random variables  $T(B_k)$ . Each variable has a variance  $m(B_k)$  (remember that T is an isometry) and is modulated in time by the harmonic function  $e^{i\omega_k t}$ , which oscillates at the frequency  $\omega_k$ .
## 6.4 Multivariable stationary processes

The theory of wide-sense stationary processes extends in a natural way to the multivariable case. This extension is dealt with here in brief summary. We consider processes with two components only because this case captures all the relevant aspects.

Let  $v_t : \Omega \to \mathbb{C}^2$  have components  $v_t^{(1)}$  and  $v_t^{(2)}$ .  $v_t$  is wide-sense stationary if

$$\mathbb{E}[v_t] = \mathbb{E}[v_0], \quad \forall t, \tag{6.56}$$

and

$$\mathbb{E}\left[\left(v_{t+\ell} - E[v_{t+\ell}]\right)\overline{\left(v_t - E[v_t]\right)^T}\right] = \mathbb{E}\left[\left(v_\ell - E[v_\ell]\right)\overline{\left(v_0 - E[v_0]\right)^T}\right], \quad \forall t, \ell, \quad (6.57)$$

where a bi-dimensional stochastic variable is identified with the vector of its components where it needs be.

Notations and terminology are the same as in the mono-variate case, so e.g. we call  $\gamma_{\ell} := \mathbb{E}\left[\left(v_{\ell} - E[v_{\ell}]\right)\overline{\left(v_{0} - E[v_{0}]\right)^{T}}\right]$  the process auto-covariance function.  $\gamma_{\ell}$  is a 2 × 2 matrix, where the diagonal elements are the auto-covariance functions of  $v_{t}^{(1)}$  and  $v_{t}^{(2)}$  and the extra-diagonal elements measure the cross-covariance between  $v_{t}^{(1)}$  and  $v_{t}^{(2)}$ . When we want to emphasize this fact, we also call  $\gamma_{\ell}^{(1,2)}$  (the (1,2) element of  $\gamma_{\ell}$ ) the cross-covariance function between  $v_{t}^{(1)}$  and  $v_{t}^{(2)}$ .

If  $v_t : \Omega \to \mathbb{C}^2$  is wide-sense stationary, so is  $\alpha v_t^{(1)} + \beta v_t^{(2)}$  for any choice of complex numbers  $\alpha$  and  $\beta$ .

The spectral theory is extended more easily to the bi-dimensional case by first introducing the spectral representation and then the spectral measure (so reversing the order adopted in the 1-dimensional case), so we follow this route.

The spectral representation is simply a componentwise concept: from the onedimensional theory,  $v_t^{(1)}$  has associated an orthogonal stochastic measure  $T^{(1)}$  such that  $v_t^{(1)} = \int_{(-\pi,\pi]} e^{i\omega t} dT^{(1)}(\omega)$ ; similarly  $v_t^{(2)} = \int_{(-\pi,\pi]} e^{i\omega t} dT^{(2)}(\omega)$ . It is worth remarking that  $T^{(1)}$  and  $T^{(2)}$  carry all information on the bi-dimensional process since  $v_t^{(1)}$  and  $v_t^{(2)}$  can be fully reconstructed from  $T^{(1)}$  and  $T^{(2)}$ . So, one can e.g. reconstruct the cross-covariance function from  $T^{(1)}$  and  $T^{(2)}$ .

In contrast, the spectral measure is a matrix concept: it is a  $2 \times 2$  matrix of the form

$$m = \begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix}.$$
 (6.58)

It needs to be so because one cannot reconstruct the cross-covariance from the autocovariance  $\gamma_{\ell}^{(1,1)}$  and  $\gamma_{\ell}^{(2,2)}$  only. The reader is invited to reflect on this point by considering the following two situations:

(i)  $v_t^{(1)}$  and  $v_t^{(2)}$  are zero-mean, unitary-variance, white and mutually uncorrelated:  $E[v_t^{(1)}\overline{v}_{\tau}^{(2)}] = 0, \forall t, \tau;$ 

(ii)  $v_t^{(1)}$  and  $v_t^{(2)}$  are zero-mean, unitary variance, white and  $v_t^{(1)} = v_t^{(2)}$ ,  $\forall t$ . Here, in both (i) and (ii)  $\gamma_{\ell}^{(1,1)} = \gamma_{\ell}^{(2,2)} = 1$  for  $\ell = 0$  and  $\gamma_{\ell}^{(1,1)} = \gamma_{\ell}^{(2,2)} = 0$  for  $\ell \neq 0$ . However,  $\gamma_{\ell}^{(1,2)}$  are different in the two cases.

Before proceeding any further in the definition of m, we need to extend our notion of measure to complex-valued signed-measures, as  $m_{12}$  and  $m_{21}$  are measures of this type.

**DEFINITION 6.11 (complex-valued signed-measure)** Let  $\mathscr{X}$  be a  $\sigma$ -algebra. A function  $m : \mathscr{X} \to [-\infty, \infty]$  such that no two sets  $A_1$  and  $A_2$  exist with  $m(A_1) = \infty$  and  $m(A_2) = -\infty$  is called a signed-measure if, for any countable collection of pairwise disjoint sets  $A_k \in \mathscr{X}$ , k = 1, 2, ..., the following property ( $\sigma$ -additivity) holds

$$m(\cup_{k=1}^{\infty}A_k) = \sum_{k=1}^{\infty}m(A_k).$$
 (6.59)

A complex-valued signed-measure m is given by  $m_1 + im_2$ , where  $m_1$  and  $m_2$  are signed-measures.

The assumption that no two sets  $B_1$  and  $B_2$  exist with  $m(B_1) = \infty$  and  $m(B_2) = -\infty$  rules out the possibility of indeterminate forms  $\infty - \infty$ . By a comparison with Definition 1.7, we see that a signed-measure is simply a measure where the positivity requirement has been relaxed.

We are now ready to define the spectral measure. For  $B \in \mathscr{B}(-\pi, \pi]$ , let:

$$m_{11}(B) = \mathbb{E}[|T^{(1)}(B)|^2], \qquad m_{12}(B) = \mathbb{E}[T^{(1)}(B)\overline{T}^{(2)}(B)], \qquad (6.60)$$
  
$$m_{21}(B) = \mathbb{E}[T^{(2)}(B)\overline{T}^{(1)}(B)], \qquad m_{22}(B) = \mathbb{E}[|T^{(2)}(B)|^2].$$

 $m_{11}$  and  $m_{22}$  are the usual spectral measures for processes  $v_t^{(1)}$  and  $v_t^{(2)}$ . From their definition, we see that  $m_{12} = \overline{m}_{21}$  and these are complex-valued signed-measures. To see this, we need to verify the  $\sigma$ -additivity property, i.e.,  $m_{12} \left( \bigcup_{k=1}^{\infty} B_k \right) = \sum_{k=1}^{\infty} m_{12} \left( B_k \right)$ , and this requires a bit of extra investigation: we shall prove that

$$\mathbb{E}[T^{(1)}(B_1)\overline{T}^{(2)}(B_2)] = 0, \quad \text{whenever } B_1 \cap B_2 = \emptyset, \tag{6.61}$$

from which the  $\sigma$ -additivity follows:

$$m_{12}(\bigcup_{k=1}^{\infty} B_k) = \mathbb{E}[T^{(1)}(\bigcup_{k=1}^{\infty} B_k)\overline{T}^{(2)}(\bigcup_{k=1}^{\infty} B_k)]$$
(6.62)

$$= \mathbb{E}\left[\left(\sum_{k=1}^{\infty} T^{(1)}(B_k)\right)\left(\sum_{k=1}^{\infty} \overline{T}^{(2)}(B_k)\right)\right]$$
(6.63)

$$= \sum_{k=1}^{\infty} \sum_{j=1}^{\infty} \mathbb{E}\left[T^{(1)}(B_k)\overline{T}^{(2)}(B_j)\right]$$
(6.64)

$$= \sum_{k=1}^{\infty} \mathbb{E}\left[T^{(1)}(B_k)\overline{T}^{(2)}(B_k)\right] \quad (use \ (6.61)) \tag{6.65}$$

$$= \sum_{k=1}^{\infty} m_{12}(B_k).$$
 (6.66)

To prove (6.61), start by considering the stationary process  $v_t = \alpha v_t^{(1)} + \beta v_t^{(2)}$  and note that the orthogonal stochastic measure *T* associated to  $v_t$  is  $\alpha T^{(1)} + \beta T^{(2)}$ . Thus,

$$0 = \mathbb{E}[T(B_1)\overline{T}(B_2)] \quad (since \ B_1 \cap B_2 = \emptyset) \tag{6.67}$$

$$= \mathbb{E}[(\alpha T^{(1)}(B_1) + \beta T^{(2)}(B_1))(\overline{\alpha}\overline{T}^{(1)}(B_2) + \overline{\beta}\overline{T}^{(2)}(B_2))]$$
(6.68)

$$= \alpha \overline{\beta} \mathbb{E}[T^{(1)}(B_1)\overline{T}^{(2)}(B_2)] + \beta \overline{\alpha} \mathbb{E}[T^{(2)}(B_1)\overline{T}^{(1)}(B_2)].$$
(6.69)

$$(since \mathbb{E}[T^{(1)}(B_1)\overline{T}^{(1)}(B_2)] = \mathbb{E}[T^{(2)}(B_1)\overline{T}^{(2)}(B_2)] = 0)$$
(6.70)

Taking  $\alpha = \beta = 1$  first, and  $\alpha = 1$ ,  $\beta = i$  then, yields

$$0 = \mathbb{E}[T^{(1)}(B_1)\overline{T}^{(2)}(B_2)] + \mathbb{E}[T^{(2)}(B_1)\overline{T}^{(1)}(B_2)]$$
(6.71)

$$0 = -i\mathbb{E}[T^{(1)}(B_1)\overline{T}^{(2)}(B_2)] + i\mathbb{E}[T^{(2)}(B_1)\overline{T}^{(1)}(B_2)], \qquad (6.72)$$

from which  $E[T^{(1)}(B_1)\overline{T}^{(2)}(B_2)] = 0$  follows, so proving (6.61).

We want to finally recall that Herglotz Theorem 6.6 extends naturally to the multi-variable case:

$$\gamma_{\ell}^{(i,j)} = \int_{(-\pi,\pi]} e^{i\omega\ell} dm_{ij}(\omega), \quad \ell = \dots, -2, -1, 0, 1, 2, \dots, \quad i, j = 1, 2.$$
(6.73)

For i = j = 1 and i = j = 2, this is the standard Herglotz theorem applied componentwise. As for  $i \neq j$ , take a partition  $B_k = (-\pi + 2\pi \frac{k-1}{N}, -\pi + 2\pi \frac{k}{N}]$ , k = 1, ..., N, of  $(-\pi, \pi]$  and let  $\sum_{k=1}^{N} e^{i\omega_k \ell} T^{(1)}(B_k)$  and  $\sum_{k=1}^{N} e^{i\omega_k 0} T^{(2)}(B_k) = \sum_{k=1}^{N} T^{(2)}(B_k)$  be  $\varepsilon$ approximations (in the  $\mathbb{L}^2$ -norm) of  $v_{\ell}^{(1)}$  and  $v_0^{(2)}$ . We have

$$\gamma_{\ell}^{(1,2)} = \mathbb{E}[v_{\ell}^{(1)}\overline{v}_{0}^{(2)}]$$
(6.74)

$$\approx \mathbb{E}\left[\left(\sum_{k=1}^{N} e^{i\omega_{k}\ell}T^{(1)}(B_{k})\right)\left(\sum_{k=1}^{N}\overline{T}^{(2)}(B_{k})\right)\right]$$
(6.75)

$$= \sum_{k=1}^{N} \sum_{j=1}^{N} e^{i\omega_{k}\ell} \mathbb{E}[T^{(1)}(B_{k})\overline{T}^{(2)}(B_{j}))]$$
(6.76)

$$= \sum_{k=1}^{N} e^{i\omega_{k}\ell} \mathbb{E}[T^{(1)}(B_{k})\overline{T}^{(2)}(B_{k}))] \quad (use \ (6.61)) \tag{6.77}$$

$$= \sum_{k=1}^{N} e^{i\omega_k \ell} m_{12}(B_k)$$
(6.78)

$$\approx \int_{(-\pi,\pi]} e^{i\omega\ell} dm_{12}(\omega), \qquad (6.79)$$

where the " $\approx$ " become "=" in the limit when  $\varepsilon \rightarrow 0$ , so proving (6.73).

## **Bibliography**

- [1] H. Bauer. *Probability theory and elements of measure theory*. Academic Press, 1981.
- [2] P. Billingsley. Probability and measure. 3rd edition, Wiley, New York, 1995.
- [3] P.R. Halmos. *Measure theory*. Van Nostrand, New York, 1950.
- [4] L.Devroye, L.Gyorfi, and G.Lugosi. *A probabilistic theory of pattern recognition*. Springer-Verlag, New York, 1996.
- [5] M. Loeve. Probability theory. Springer-Verlag, New York, 1977.
- [6] W. Rudin. Real and complex analysis. McGraw-Hill, New York, 1966.
- [7] A.N. Shiryaev. Probability. Springer, 2nd edition, 1996.