

Randal Douc and Sylvain Le Corff

# Bayesian Learning for partially observed dynamical systems

September 23, 2019



# Contents

<b>1</b>	<b>Preliminaries</b>	3
1.1	Refresher on Martingales	4
1.2	Some usual distributions	5
1.3	Abbreviations	5
<b>2</b>	<b>Introduction to Markov Chains</b>	7
2.1	Kernels	8
2.2	Homogeneous Markov chain	10
2.3	Canonical representation	12
2.4	Invariant measures	12
2.5	Observation-driven models	14
2.6	Iterated random functions	16
<b>3</b>	<b>Inference for Markovian Models</b>	29
3.1	Likelihood inference	29
3.2	Consistency and asymptotic normality of the MLE	33
3.3	Observation-driven models	41
3.4	Bayesian inference	48
3.5	Some proofs	55
	Exercises	58
<b>4</b>	<b>MCMC methods</b>	65
4.1	Metropolis-Hastings algorithm	66
4.2	Ordering the asymptotic variances	76
<b>5</b>	<b>Ergodic theory for Markov chains</b>	79
5.1	Dynamical systems	80
5.2	Markov chains ergodicity	84
5.3	Exercises	89
5.4	Bibliographical notes	93
<b>6</b>	<b>Pseudo Marginal Monte Carlo methods and applications</b>	95
<b>7</b>	<b>Hamiltonian Monte Carlo methods</b>	97
7.1	MH with deterministic moves	97
7.2	Hamiltonian dynamics	98
7.3	The leapfrog integrator	99
<b>8</b>	<b>Exercises</b>	101
	Exercises	101

<b>Appendix A    Technical results</b> .....	109
A.1    Limit theorems for triangular arrays .....	110
<b>References</b> .....	115

# Plan of action

This is a 21 hours course schedule. The plan of action is the following.

1. Markov chains: kernel, invariant measures. Examples including Observation-Driven Models: lecture: 1H45, tutorials: 1H45. Randal Cours. Sylvain TD
2. Bayesian inference, asymptotic properties of the MLE for Markov chains : lecture: 1H45, tutorials: 1H45. Sylvain Cours et Td.
3. Markov Chains Monte Carlo (MCMC) algorithms: lecture: 1H45, computer sessions: 1H45. Randal Cours, Td.
4. Some properties of the MCMC algorithms: lecture: 1H45, tutorials: 1H45. Randal Cours.
5. Pseudo marginal MCMC and applications: lecture: 1H45, computer sessions: 1H45. Sylvain Cours, td.
6. Hamiltonian Monte Carlo algorithms: lecture: 1H30, tutorials and computer sessions: 1H45 Sylvain Cours, td.



# Chapter 1

## Preliminaries

### Contents

1.1	Refresher on Martingales .....	4
1.2	Some usual distributions .....	5
1.3	Abbreviations .....	5

These lecture notes are built from various sources:

- (i) The monograph *Nonlinear Time Series: Theory, Methods, and Applications with R Examples*, by R. Douc, E. Moulines and D. Stoffer. Chapman and Hall editors.
- (ii) *Handbook of Markov Chain Monte Carlo* by Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. Chapman and Hall editors.

Do not hesitate to point out the errors or typos that still remain and to propose any improvements on the content of this course. Email contact:

- [randal.douc@telecom-sudparis.eu](mailto:randal.douc@telecom-sudparis.eu)
- [sylvain.lecorff@telecom-sudparis.eu](mailto:sylvain.lecorff@telecom-sudparis.eu)

We start with **some notation**. In what follows,

- i.i.d means independent and identically distributed.
- r.v. means random variables.
- for  $r, s \in \mathbb{N}$  such that  $r \leq s$ , we write  $[r : s] = \{r, r+1, \dots, s\}$ ,
- $X \perp\!\!\!\perp Y$  means  $X$  and  $Y$  are independent random variables,
- $X \stackrel{\mathcal{L}}{=} Y$  means  $X$  and  $Y$  have the same law.
- Let  $(Z, \mathcal{Z})$  and  $(X, \mathcal{X})$  two measurable spaces. Assume that for all  $w \in Z$ , the two random vectors  $X$  and  $Y(w)$  take values in  $X$  and assume the existence of a third random variable  $W$  taking values in  $Z$ , then the notation:  $X|_{W=w} \stackrel{\mathcal{L}}{=} Y(w)$  means that for all  $A \in \mathcal{X}$ ,

$$\mathbb{P}(X \in A | W) |_{W=w} = \mathbb{P}(Y(w) \in A)$$

In words, the distribution of  $X$  conditionally on  $W$  taken on  $W = w$  is the same as the unconditional distribution of  $Y(w)$ .

- $\liminf_n a_n = \lim_{n \rightarrow \infty} (\inf_{k \geq n} a_k)$  and similarly,  $\limsup_n a_n = \lim_{n \rightarrow \infty} (\sup_{k \geq n} a_k)$ . Moreover,  $\lim_n a_n$  exists if and only if  $\liminf_n a_n = \limsup_n a_n$ .
- for any  $a \in \mathbb{R}$ ,  $a^+ = \max(a, 0)$  and  $a^- = \max(-a, 0) = -\min(a, 0)$  and we have  $|a| = a^+ + a^-$  and  $a = a^+ - a^-$ .

Moreover, the following notions of convergence for random variables is used throughout these lecture notes.

- $X_n \xrightarrow{w} X$  means *convergence in distribution* (or “convergence en loi” in French). It is equivalent to any of the following statements.

- (a) for all bounded continuous functions  $h$ , we have  $\lim_n \mathbb{E}[h(X_n)] = \mathbb{E}[h(X)]$ .
- (b) for all  $A \in \mathcal{B}(\mathbb{R})$  such that  $\mathbb{P}(X \in \partial A) = 0$ , we have  $\lim_n \mathbb{P}(X_n \in A) = \mathbb{P}(X \in A)$ .
- (c) for all  $x \in \mathbb{R}$  such that  $\mathbb{P}(X = x) = 0$ , we have  $\lim_n \mathbb{P}(X_n \leq x) = \mathbb{P}(X \leq x)$ .
- (d) for all  $u \in \mathbb{R}$ , we have  $\lim_n \mathbb{E}[e^{iuX_n}] = \mathbb{E}[e^{iuX}]$

By abuse of terminology, we may also say that  $X_n$  **weakly converges to**  $X$  instead of saying the distribution of  $X_n$  converges weakly to the distribution of  $X$ . An equivalent formulation is  $X_n \xrightarrow{\mathcal{L}_P} X$ .

- $X_n \xrightarrow{\mathbb{P}\text{-prob}} X$  means *convergence in probability*: for all  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0.$$

- $X_n \xrightarrow{\mathbb{P}\text{-a.s.}} X$  means *almost sure convergence*:

$$\mathbb{P}(\lim_{n \rightarrow \infty} X_n = X) = 1.$$

Almost sure convergence implies convergence in probability. We recall the following properties.

- (i) If  $X_n \xrightarrow{w} X$  then for all continuous functions  $f$ ,  $f(X_n) \xrightarrow{w} f(X)$ . Note that this property holds, when  $f$  is continuous (and not necessarily bounded), for example  $f(u) = u^2$  so that  $X_n^2 \xrightarrow{w} X^2$ .
- (ii) **The Slutsky Lemma** If  $X_n \xrightarrow{\mathbb{P}\text{-prob}} c$  where  $c$  is a constant and if  $Y_n \xrightarrow{w} Y$ , then  $(X_n, Y_n) \xrightarrow{w} (c, Y)$  that is for all continuous functions  $f$ ,  $f(X_n, Y_n) \xrightarrow{w} f(c, Y)$ .
- (iii)  $X \sim N(0, 1)$  iff for all  $u \geq 0$ ,  $\mathbb{E}[e^{iuX}] = e^{-u^2/2}$ . Moreover,  $X \sim N(\mu, \sigma^2)$  iff for all  $u \geq 0$ ,  $\mathbb{E}[e^{iuX}] = e^{-u^2 \text{Var}(X)/2 + iu\mathbb{E}(X)}$  and in that case,  $\sigma^2 = \text{Var}(X)$  and  $\mu = \mathbb{E}(X)$ .

## 1.1 Refresher on Martingales

- “ $\{\mathcal{F}_t, t \in \mathbb{N}\}$  is a filtration” means that  $(\mathcal{F}_t)$  is an increasing sequence of  $\sigma$ -fields, that is  $\mathcal{F}_t$  is a  $\sigma$ -field and  $\mathcal{F}_t \subset \mathcal{F}_{t+1}$  for all  $t \in \mathbb{N}$ .
- “ $(\Omega, \mathcal{F}, \{\mathcal{F}_t, t \in \mathbb{N}\}, \mathbb{P})$  is a filtered probability space” means that  $\mathcal{F}$  is a  $\sigma$ -field on  $\Omega$ ,  $\mathbb{P}$  is a probability on  $(\Omega, \mathcal{F})$  and  $\{\mathcal{F}_t, t \in \mathbb{N}\}$  is a filtration such that  $\mathcal{F}_t \subset \mathcal{F}$  for all  $t \in \mathbb{N}$ .

Moreover, the sequence  $(M_t)_{t \in \mathbb{N}}$  is a martingale wrt the filtration  $\{\mathcal{F}_t, t \in \mathbb{N}\}$ , if for all  $t \geq 0$ ,  $M_t$  is a random variable on  $(\Omega, \mathcal{F}, \mathbb{P})$  and

$$\mathbb{E}[M_{t+1} | \mathcal{F}_t] = M_t.$$

The following property on martingales will be important in this course.

**Theorem 1.1.** *Let  $(M_t)_{t \in \mathbb{N}}$  be a martingale wrt the filtration  $\{\mathcal{F}_t, t \in \mathbb{N}\}$  and denote  $\Delta_t = M_t - M_{t-1}$  for  $t \geq 1$ . Assume that  $\{\Delta_k, k \in \mathbb{N}\}$  is a strict-sense stationary and ergodic process such that  $\mathbb{E}[M_1^2]$  is finite. Then,*

$$n^{-1/2} M_n \xrightarrow{\mathcal{L}_P} N(0, \mathbb{E}[\Delta_1^2]).$$

Another interesting property is actually a byproduct of the martingale theory since  $M_n = \mathbb{E}[\mathbb{1}_A | \mathcal{F}_n]$  is a nonnegative martingale.

**Lemma 1.2** *Let  $A \in \mathcal{F} = \sigma(\cup_{k=0}^\infty \mathcal{F}_k)$ . Then,*

$$\lim_{n \rightarrow \infty} \mathbb{E}[\mathbb{1}_A | \mathcal{F}_n] = \mathbb{1}_A \quad \mathbb{P} - \text{a.s.}$$



## 1.2 Some usual distributions

Name	Acronym	Parameter	density function: $f_X(x)$	cdf: $F_X(x) = \int_{-\infty}^x f_X(u)du$	Other properties
Gaussian	$N(\mu, \sigma^2)$	$(\mu, \sigma^2)$	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)}$	No explicit expression	$\mathbb{E}[X] = \mu, \mathbb{V}\text{ar}X = \sigma^2$
Exponential	$\exp(\lambda)$	$\lambda > 0$	$\lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}_+}(x)$	$(1 - e^{-\lambda x}) \mathbb{1}_{\mathbb{R}_+}(x)$	$\mathbb{E}[X] = 1/\lambda, \mathbb{V}\text{ar}X = 1/\lambda^2$
Gamma	$\Gamma(k, \theta)$	$(k, \theta) \in (\mathbb{R}_+^*)^2$	$\frac{x^{k-1} e^{-x/\theta}}{\Gamma(k)\theta^k}$	$\frac{\Gamma_{x/\theta}(k)}{\Gamma(k)}$	

In the above description,

- (i) if  $X_i \sim \Gamma(k_i, \theta)$  and  $(X_i)$  are independent, then  $\sum_{i=1}^n X_i \sim \Gamma(\sum_{i=1}^n k_i, \theta)$ .  
(ii)

$$\Gamma(k) = \begin{cases} \int_0^\infty t^{k-1} e^{-t} dt & \text{if } k \in \mathbb{R}_+^* \\ k! & \text{if } k \in \mathbb{N}. \end{cases} \quad (\blacktriangleright \text{GAMMA FUNCTION})$$

$$\Gamma_x(k) = \int_0^x t^{k-1} e^{-t} dt \quad (\blacktriangleright \text{INCOMPLETE GAMMA FUNCTION})$$

## 1.3 Abbreviations

$\mathbb{F}(X)$	The set of measurable functions on $(X, \mathcal{X})$
$\mathbb{F}_+(X)$	The set of nonnegative measurable functions on $(X, \mathcal{X})$
$\mathbb{M}_1(\mathcal{X})$	The set of probability measures on $(X, \mathcal{X})$
$\mathbb{M}_+(\mathcal{X})$	The set of (nonnegative) measures on $(X, \mathcal{X})$



Chapter

2

# Introduction to Markov Chains

## Contents

2.1	Kernels .....	8
2.2	Homogeneous Markov chain .....	10
2.3	Canonical representation .....	12
2.4	Invariant measures .....	12
2.5	Observation-driven models .....	14
2.6	Iterated random functions .....	16
2.6.1	Strict stationarity .....	16
2.6.2	Weak stationarity .....	19
2.6.3	Iterated random functions .....	21

There are two rather different approaches to constructing a Markov chain. The first is through the transition laws of the chain. The second is through the use of iterating functions. This is particularly suited to the use of Markov chains in a time series context, where we wish to construct sample paths of the chain in an explicit manner rather than work with the distributions of the chain. This can also be carried out in a great degree of generality.

For both of these approaches we will discuss the existence of invariant measures (sometimes called stationary distributions) for the chain. This is the key property for most models because it implies (and is in some sense equivalent to) the existence of strict-sense stationary solutions. More precisely, when a Markov chain admits an invariant distribution, then we may construct strictly stationary versions of Markov chains; these stationary versions play a key role when considering the inference of Markov chains.

We will illustrate these constructions with several important time series models. It is worthwhile to note that, while many models in the nonlinear time series context might be viewed as (most of the time straightforward) extensions of ARMA models, finding necessary and sufficient conditions for the existence of an invariant distribution is a difficult task that will be covered in some generality in subsequent chapters.

Some of the material in this chapter will seem, to those interested in applications, to be both technical and exotic. However, we use the formal structures to define, in a rigorous manner, different time-series models. We encourage the reader to go through this material, perhaps skipping the proofs in a first pass.

A Markov chain is a discrete time stochastic process  $X = \{X_t, t \in \mathbb{N}\}$ , i.e., a countable collection of random variables defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . In this definition,  $t$  is thought of as a time index and the set of times is taken by convention as  $\mathbb{N}$ .

The chain  $X$  evolves on a *state space*  $X$ . Although we let  $X$  be a general set (not assuming a priori any topological structure), in practice, it is most likely that we will be considering the set of real numbers or, more generally, the  $d$ -dimensional Euclidean space, so that  $X = \mathbb{R}^d$ , or some countable or uncountable subset of  $\mathbb{R}^d$ . In order to define probabilities, we use  $\mathcal{X}$  to denote a countably generated  $\sigma$ -field on  $X$ . When  $X = \mathbb{R}^d$ , then  $\mathcal{X}$  will be taken as the Borel  $\sigma$ -field  $\mathcal{X} = \mathcal{B}(X)$ . When  $X$  is countable,  $\mathcal{X}$  contains all the subsets, that is  $\mathcal{X} = \mathcal{P}(X)$ , the set of all parts of  $X$ . All random variables are assumed measurable individually with respect to  $\mathcal{B}(X)$ , and we shall in general denote elements of  $X$  by letters  $x, y, z, \dots$  and elements of  $\mathcal{X}$  by  $A, B, C$ .

The distinguishing feature of a Markov chain as opposed to an arbitrary stochastic process is that the future, given the present, is forgetful of the past. This is reflected in the conditional distributions of the future random variable  $X_{t+1}$  given its present,  $X_t$ , and its past  $\{X_0, \dots, X_{t-1}\}$ .

**Definition 2.1 (Markov chain).** Let  $(\Omega, \mathcal{F}, \{\mathcal{F}_t, t \in \mathbb{N}\}, \mathbb{P})$  be a filtered probability space. An adapted stochastic process  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is called a Markov chain of *order  $m$*  (or with *memory  $m$* ) if, for all  $t \geq m-1$  and  $A \in \mathcal{X}$ ,  $\mathbb{P}[X_{t+1} \in A | \mathcal{F}_t] = \mathbb{P}[X_{t+1} \in A | X_t, \dots, X_{t-m+1}]$ ,  $\mathbb{P}$ -a.s.

When the order is one, we simply say that  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is a Markov chain. This definition is equivalent to assuming that for all  $f \in \mathbb{F}_+(\mathbf{X})$ ,

$$\mathbb{E}[f(X_t) | \mathcal{F}_{t-1}] = \mathbb{E}[f(X_t) | X_{t-1}], \quad \mathbb{P} - \text{a.s.}$$

Note that a Markov chain  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is always a Markov chain with respect to the *natural filtration*  $\{\mathcal{F}_t^X, t \in \mathbb{N}\}$ , where  $\mathcal{F}_t^X = \sigma(X_0, \dots, X_t)$ .

## 2.1 Kernels

The construction of Markov chains is based on the definition of kernels.

**Definition 2.2.** Let  $(X, \mathcal{X})$  and  $(Y, \mathcal{Y})$  be two measurable spaces.

- A *kernel* is a mapping  $N$  from  $X \times \mathcal{Y}$  into  $[0, \infty]$  satisfying the following conditions:
  - (i) for every  $x \in X$ , the mapping  $N(x, \cdot) : A \mapsto N(x, A)$  is a measure on  $\mathcal{Y}$ ,
  - (ii) for every  $A \in \mathcal{Y}$  the mapping  $N(\cdot, A) : x \mapsto N(x, A)$  is a measurable function from  $(X, \mathcal{X})$  to  $[0, \infty]$ .
- The kernel  $N$  is said to be *finite* if  $N(x, Y) < \infty$  for all  $x \in X$ .
- The kernel  $N$  is said to be *bounded* if  $\sup_{x \in X} N(x, Y) < \infty$ .
- A kernel  $N$  is said to be *Markovian* if  $N(x, Y) = 1$ , for all  $x \in X$ .

**Example 2.3 (Discrete state-space kernel)** Assume that  $X$  and  $Y$  are countable sets. Each element  $x \in X$  is then called a *state*. A kernel  $N$  on  $X \times \mathcal{P}(Y)$ , where  $\mathcal{P}(Y)$  is the set of all parts of  $Y$ , is specified by a (possibly infinite) matrix  $N = (N(x, y) : x, y \in X \times Y)$  with nonnegative entries. Each row  $(N(x, y) : y \in Y)$  is a measure on  $(Y, \mathcal{P}(Y))$  defined by

$$N(x, A) = \sum_{y \in A} N(x, y),$$

for  $A \subset Y$ . The matrix  $N$  is said to be *Markovian* if every row  $(N(x, y) : y \in Y)$  is a probability measure on  $(Y, \mathcal{P}(Y))$ , i.e.  $\sum_{y \in Y} N(x, y) = 1$  for all  $x \in X$ .

**Example 2.4 (Kernel density)** Let  $\lambda \in \mathbb{M}_1(\mathcal{Y})$  and  $n : X \times Y \rightarrow \mathbb{R}_+$  be a nonnegative function, measurable with respect to the product  $\sigma$ -algebra  $\mathcal{X} \otimes \mathcal{Y}$ . Then, the application  $N$  defined on  $X \times \mathcal{Y}$  by

$$N(x, A) = \int_A n(x, y) \lambda(dy),$$

is a kernel. The function  $n$  is called the *density of the kernel  $N$*  with respect to the measure  $\lambda$ . The kernel  $N$  is *Markovian* if and only if  $\int_Y n(x, y) \lambda(dy) = 1$  for all  $x \in X$ .

Let  $N$  be a kernel and  $f \in \mathbb{F}_+(\mathcal{Y})$ . A function  $Nf : \mathcal{X} \rightarrow \mathbb{R}^+$  is defined by setting, for any  $x \in \mathcal{X}$ ,

$$Nf : x \mapsto \int_{\mathcal{Y}} N(x, dy) f(y) .$$

**Proposition 2.5** *Let  $N$  be a kernel on  $\mathcal{X} \times \mathcal{Y}$ . For any  $f \in \mathbb{F}_+(\mathcal{Y})$ ,  $Nf \in \mathbb{F}_+(\mathcal{X})$ .*

With a slight abuse of notation, we will use the same symbol  $N$  for the kernel and the associated operator  $N : \mathbb{F}_+(\mathcal{Y}) \rightarrow \mathbb{F}_+(\mathcal{X})$ ,  $f \mapsto Nf$ . By defining  $Nf = Nf^+ - Nf^-$ , we may extend the application  $f \mapsto Nf$  to all functions  $f$  of  $\mathbb{F}(\mathcal{Y})$  such that  $Nf^+$  and  $Nf^-$  are not both infinite. We will sometimes write  $N(x, f)$  for  $Nf(x)$  and make use of the notations  $N(x, \mathbb{1}_A)$  and  $N\mathbb{1}_A(x)$  for  $N(x, A)$ .

Let  $\mu$  be a positive measure on  $(\mathcal{X}, \mathcal{X})$  and for  $A \in \mathcal{Y}$ , define

$$\mu N(A) = \int_{\mathcal{X}} \mu(dx) N(x, A) .$$

**Proposition 2.6** *Let  $N$  be a kernel on  $\mathcal{X} \times \mathcal{Y}$  and  $\mu \in \mathbb{M}_+(\mathcal{X})$ . Then  $\mu N \in \mathbb{M}_+(\mathcal{Y})$ .*

Let  $(\mathcal{X}, \mathcal{X})$ ,  $(\mathcal{Y}, \mathcal{Y})$  and  $(\mathcal{Z}, \mathcal{Z})$  be measurable spaces and let  $M$  and  $N$  be two kernels on  $\mathcal{X} \times \mathcal{Y}$  and  $\mathcal{Y} \times \mathcal{Z}$ . For any  $A \in \mathcal{Z}$ ,  $y \mapsto N(y, A)$  is a measurable function, and by Proposition 2.5,  $x \mapsto \int_{\mathcal{Y}} M(x, dy) N(y, A)$  is a measurable function. For any  $x \in \mathcal{X}$ ,  $M(x, \cdot)$  is a measure on  $(\mathcal{Y}, \mathcal{Y})$  and by Proposition 2.6,  $A \mapsto \int_{\mathcal{Y}} M(x, dy) N(y, A)$  is a measure on  $(\mathcal{Z}, \mathcal{Z})$ . Hence, the function  $MN : (x, A) \mapsto \int_{\mathcal{Y}} M(x, dy) N(y, A)$  is a kernel on  $\mathcal{X} \times \mathcal{Z}$ .

The *composition* or *product* of the kernels  $M$  on  $\mathcal{X} \times \mathcal{Y}$  and  $N$  on  $\mathcal{Y} \times \mathcal{Z}$  is the kernel  $MN$  defined for  $x \in \mathcal{X}$  and  $A \in \mathcal{Z}$  by

$$MN(x, A) = \int_{\mathcal{Y}} M(x, dy) N(y, A) . \quad (2.1)$$

Since  $MN$  is a kernel on  $\mathcal{X} \times \mathcal{Z}$ , for any  $f \in \mathbb{F}_+(\mathcal{Z})$ , we may define the function  $MNf : x \mapsto MNf(x)$ , which by Proposition 2.5 belongs to  $\mathbb{F}_+(\mathcal{X})$ . On the other hand,  $Nf : y \mapsto Nf(y)$  is a function belonging to  $\mathbb{F}_+(\mathcal{Y})$ , and since  $M$  is a kernel on  $\mathcal{X} \times \mathcal{Y}$ , we may consider the function  $x \mapsto M[Nf](x)$ . As shown in the following proposition, these two quantities coincide.

**Proposition 2.7** *Let  $M$  be a kernel on  $\mathcal{X} \times \mathcal{Y}$  and  $N$  be a kernel on  $\mathcal{Y} \times \mathcal{Z}$ . Then, for each  $x \in \mathcal{X}$ , and  $f \in \mathbb{F}_+(\mathcal{Z})$ ,*

$$MNf(x) = M[Nf](x) . \quad (2.2)$$

For  $x \in \mathcal{X}$  and  $A \in \mathcal{Z}$ , we set  $N^0(x, A) = \delta_x(A)$  and for  $n \geq 1$ , we define inductively  $N^n$  by

$$N^n(x, A) = \int_{\mathcal{X}} N(x, dy) N^{n-1}(y, A) . \quad (2.3)$$

We finally define the tensor products of kernels.

**Proposition 2.8** *Let  $M$  be a kernel on  $\mathcal{X} \times \mathcal{Y}$  and  $N$  be a kernel on  $\mathcal{Y} \times \mathcal{Z}$ . Then, there exists a kernel  $M \otimes N$  on  $\mathcal{X} \times (\mathcal{Y} \otimes \mathcal{Z})$ , called the tensor product of  $M$  and  $N$ , such that, for all  $f \in \mathbb{F}_+(\mathcal{Y} \times \mathcal{Z}, \mathcal{Y} \otimes \mathcal{Z})$ ,*

$$M \otimes Nf(x) = \int_{\mathcal{Y}} M(x, dy) \int_{\mathcal{Z}} f(y, z) N(y, dz) . \quad (2.4)$$

For  $n \geq 1$ , the  $n$ -th tensorial power  $P^{\otimes n}$  of a kernel  $P$  on  $\mathcal{X} \times \mathcal{Y}$  is the kernel on  $(\mathcal{X}, \mathcal{X}^{\otimes n})$  defined by

$$P^{\otimes n}f(x) = \int_{\mathcal{X}^n} f(x_1, \dots, x_n) P(x, dx_1) P(x_1, dx_2) \cdots P(x_{n-1}, dx_n) .$$

## 2.2 Homogeneous Markov chain

**Definition 2.9 (Homogeneous Markov chain).** Let  $(\Omega, \mathcal{F}, \{\mathcal{F}_t, t \in \mathbb{N}\}, \mathbb{P})$  be a filtered probability space. Let  $m \geq 1$  be an integer and  $P$  a Markov kernel on  $X^m \times \mathcal{X}$ . An adapted stochastic process  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is called a *homogeneous Markov chain of order  $m$*  with Markov kernel  $P$  and initial distribution  $\nu \in \mathbb{M}_1(\mathcal{X}^{\otimes m})$  if

- (i) for all  $A \in \mathcal{X}^{\otimes m}$ ,  $\mathbb{P}((X_0, \dots, X_{m-1}) \in A) = \nu(A)$ ,
- (ii) for all  $t \geq m-1$   $\mathbb{P}[X_{t+1} \in A | \mathcal{F}_t] = P(X_t, \dots, X_{t-m+1}, A)$ ,  $\mathbb{P}$ -a.s.

By construction, for any function  $f_0 \in \mathbb{F}_+(X)$ ,  $\mathbb{E}[f_0(X_0)] = \int \nu(dx_0) f_0(x_0) = \nu(f_0)$ . For all functions  $f_0, f_1 \in \mathbb{F}_+(X)$ , we may write similarly

$$\begin{aligned} \mathbb{E}[f_0(X_0)f_1(X_1)] &= \mathbb{E}[f_0(X_0)\mathbb{E}[f_1(X_1) | \mathcal{F}_0]] = \mathbb{E}[f_0(X_0)Pf_1(X_0)] \\ &= \iint \nu(dx_0)P(x_0, dx_1)f_0(x_0)f_1(x_1) = \nu \otimes P(f_0 \otimes f_1), \end{aligned}$$

where  $f_0 \otimes f_1(x, x') = f_0(x)f_1(x')$ . Assume that for some  $t \geq 1$ , and for all functions  $f_0, \dots, f_t \in \mathbb{F}_+(X)$ ,

$$\begin{aligned} \mathbb{E}[f_0(X_0)f_1(X_1)\dots f_t(X_t)] &= \int \nu(dx_0) \prod_{s=1}^t P(x_{s-1}, dx_s) \prod_{s=0}^t f_s(x_s) \\ &= \nu \otimes P^{\otimes t}(f_0 \otimes f_1 \otimes \dots \otimes f_t). \end{aligned} \quad (2.5)$$

Then, for any  $f_{t+1} \in \mathbb{F}_+(X)$ ,

$$\begin{aligned} \mathbb{E}[f_0(X_0)\dots f_t(X_t)f_{t+1}(X_{t+1})] &= \mathbb{E}[f_0(X_0)\dots f_t(X_t)\mathbb{E}[f_{t+1}(X_{t+1}) | \mathcal{F}_t]] \\ &= \mathbb{E}[f_0(X_0)\dots f_t(X_t)Pf_{t+1}(X_{t+1})] = \nu \otimes P^{\otimes(t+1)}(f_0 \otimes f_1 \otimes \dots \otimes f_{t+1}), \end{aligned}$$

showing that (2.5) is still true with  $t$  replaced by  $t+1$ , and hence is satisfied for all  $t \in \mathbb{N}$ . Since the cylinder sets generate the product  $\sigma$ -algebra  $\mathcal{X}^{\otimes \mathbb{N}}$ , the initial distribution  $\nu$  on  $(X, \mathcal{X})$  and the Markov kernel  $P$  allows to define a unique distribution on  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$ .

Conversely, assume that for all  $t \in \mathbb{N}$ , (2.5) is satisfied. Then, using the Tonelli-Fubini theorem, we may write

$$\begin{aligned} \mathbb{E}[f_0(X_0)\dots f_{t-1}(X_{t-1})f_t(X_t)] &= \int \nu(dx_0) \prod_{s=1}^{t-1} P(x_{s-1}, dx_s) \prod_{s=0}^{t-1} f_s(x_s) \int P(x_{t-1}, dx_t) f_t(x_t) \\ &= \mathbb{E}[f_0(X_0)\dots f_{t-1}(X_{t-1})Pf_t(X_{t-1})], \end{aligned}$$

showing that for any  $\mathcal{F}_{t-1}$ -measurable random variable  $Y = f_0(X_0)\dots f_{t-1}(X_{t-1})$ ,  $\mathbb{E}[Yf_t(X_t)] = \mathbb{E}[YPf_t(X_{t-1})]$ , which implies that  $\mathbb{E}[f_t(X_t) | \mathcal{F}_{t-1}] = Pf_t(X_{t-1})$ . Hence, an adapted stochastic process satisfying (2.5) is an homogeneous Markov chain with initial distribution  $\nu$  and Markov kernel  $P$ . We summarize these conclusions in the following theorem.

**Theorem 2.10.** Let  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  be a Markov chain on  $(X, \mathcal{X})$  with initial distribution  $\nu$  and Markov kernel  $P$ . For any  $f \in \mathbb{F}_b(X^{t+1}, \mathcal{X}^{\otimes(t+1)})$ , and  $t \in \mathbb{N}$ , we have

$$\mathbb{E}[f(X_0, \dots, X_t)] = \nu \otimes P^{\otimes t}(f). \quad (2.6)$$

Conversely, let  $\{X_t, t \in \mathbb{N}\}$  be a stochastic process on  $(X, \mathcal{X})$  satisfying (2.6) for some probability  $\nu$  and a Markov kernel  $P$  on  $X \times \mathcal{X}$ . Then,  $\{X_t, t \in \mathbb{N}\}$  is a Markov chain with initial distribution  $\nu$  and transition probability  $P$ .

PROOF. The essence of the proof is presented above; technical details are filled in Exercise 8.10 and Exercise 8.11. ■

**Example 2.11 (Autoregressive process of order 1)** Consider the AR(1) process

$$X_{t+1} = \phi X_t + Z_{t+1}, \quad t \in \mathbb{N}, \quad (2.7)$$

where  $\{Z_t, t \in \mathbb{N}^*\}$  is a sequence of zero-mean i.i.d. random variables independent from  $X_0$ . For  $t \geq 0$ , define  $\mathcal{F}_t = \sigma(X_0, Z_s, s \leq t)$ . The process  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  given by (2.7) is a Markov chain. For any  $x \in X$  and  $A \in \mathcal{B}(\mathbb{R})$ , the Markov kernel of this chain is given by  $P(x, A) = \mathbb{E}[\mathbb{1}_A(\phi x + Z_0)] = \mu(A - \phi x)$ , where  $\mu$  is the distribution of  $Z_0$ . If the law of  $Z_0$  has a density  $q$ , with respect to Lebesgue's measure on  $\mathbb{R}$ , then the Markov kernel also has a density  $p(x, y) = q(y - \phi x)$ , i.e.

$$P(x, A) = \int_A q(y - \phi x) \text{Leb}(dy).$$

**Example 2.12 (ARCH(1) process)** An autoregressive conditional heteroscedastic (ARCH) model of order 1 is defined as

$$X_t = \sigma_t Z_t, \quad \sigma_t^2 = \alpha_0 + \alpha_1 X_{t-1}^2, \quad t \geq 1, \quad (2.8)$$

where the coefficients  $\alpha_0, \alpha_1$  are positive and  $\{Z_t, t \in \mathbb{N}\}$  is an i.i.d. sequence such that  $\mathbb{E}[Z_0] = 0$ ,  $\mathbb{E}[Z_0^2] = 1$ , and  $\{Z_t, t \in \mathbb{N}\}$  is independent of  $X_0$ . Equation (2.8) may be equivalently rewritten as

$$X_t = \sqrt{\alpha_0 + \alpha_1 X_{t-1}^2} Z_t,$$

which shows that  $\{X_t, t \in \mathbb{N}\}$  is an homogeneous Markov chain. The Markov kernel of this chain is given by  $P(x, A) = \mathbb{P}(\sqrt{\alpha_0 + \alpha_1 x^2} Z \in A)$ , for  $x \in \mathbb{R}$  and  $A \in \mathcal{B}(\mathbb{R})$ . If the distribution of  $Z_0$  has a density  $q$  with respect to the Lebesgue measure on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , then the Markov kernel  $P$  has a density  $p$  given by

$$p(x, x') = \frac{1}{\sqrt{\alpha_0 + \alpha_1 x^2}} q\left(\frac{x'}{\sqrt{\alpha_0 + \alpha_1 x^2}}\right).$$

The square volatility  $\{\sigma_t^2, t \geq 1\}$  satisfies the recursion: for  $t \geq 2$ ,

$$\sigma_t^2 = \alpha_0 + \alpha_1 Z_{t-1}^2 \sigma_{t-1}^2$$

and  $\sigma_1^2 = \alpha_0 + \alpha_1 X_0^2$ . This is again a Markov chain with state space  $\mathbb{R}^+$  and Markov kernel  $Q$  defined by  $Q(s, A) = \mathbb{P}(\alpha_0 + \alpha_1 s Z^2 \in A)$  for all  $s \in \mathbb{R}^+$  and  $A \in \mathcal{B}(\mathbb{R}^+)$ .

**Example 2.13 (A simple Markov bilinear process)** Consider the following bilinear process

$$X_t = aX_{t-1} + bZ_t X_{t-1} + Z_t,$$

where  $\{Z_t, t \in \mathbb{N}\}$  is an i.i.d. sequence independent from  $X_0$ . The Markov kernel of the chain  $\{X_t, t \in \mathbb{N}\}$  is given by  $P(x, A) = \mathbb{P}(ax + (1 + bx)Z \in A)$  for  $x \in \mathbb{R}$  and  $A \in \mathcal{B}(\mathbb{R})$ . If the distribution of  $Z_0$  has a density  $q$  with respect to the Lebesgue measure and if  $x \neq -1/b$ , then the Markov kernel  $P$  has a density, denoted  $p$ , given by

$$p(x, x') = \frac{1}{|1 + bx|} q\left(\frac{x' - ax}{1 + bx}\right).$$

If  $x = -1/b$ , then  $P(-1/b, \cdot) = \delta_{-a/b}$ .

### 2.3 Canonical representation

In this section, we show that, given an initial distribution  $\nu \in \mathbb{M}_1(\mathcal{X})$  and a Markov kernel  $P$  on  $X \times \mathcal{X}$ , we can construct a filtered probability space and a Markov chain with initial distribution  $\nu$  and transition kernel  $P$ . A first construction has been given by Theorem 2.37, under some topological restriction. The following construction imposes no additional assumption on  $(X, \mathcal{X})$ . Let  $X^{\mathbb{N}}$  be the set of  $X$ -valued sequences  $w = (w_0, \dots, w_n, \dots)$  endowed with the product  $\sigma$ -field  $\mathcal{X}^{\otimes \mathbb{N}}$ . We define the coordinate process  $\{X_t, t \in \mathbb{N}\}$  by

$$X_t(w) = w_t. \quad (2.9)$$

The natural filtration of the coordinate process is denoted by  $\{\mathcal{F}_t^X, t \in \mathbb{N}\}$ ; by construction,  $\bigvee_{n \geq 0} \mathcal{F}_n^X = \mathcal{X}^{\otimes \mathbb{N}}$ .

**Theorem 2.14.** *Let  $P$  be a Markov kernel  $P$  on  $X \times \mathcal{X}$  and  $\nu \in \mathbb{M}_1(\mathcal{X})$ . Then, there exists a unique probability measure  $\mathbb{P}_\nu$  on  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  such that the coordinate process is a Markov chain with initial distribution  $\nu$  and Markov kernel (or transition probability)  $P$ .*

PROOF. Set for  $t \in \mathbb{N}$  and  $f \in \mathbb{F}_+(\mathcal{X}^{t+1}, \mathcal{X}^{\otimes(t+1)})$ ,

$$\mu_t(f) = \nu \otimes P^{\otimes t}(f) = \int \nu(dx_0) \int P(x_0, dx_1) \dots \int P(x_{t-1}, dx_t) f(x_0, x_1, \dots, x_t).$$

By Proposition 2.8,  $\mu_t \in \mathbb{M}_1(\mathcal{X}^{\otimes(t+1)})$  and the family of finite dimensional distributions  $\{\mu_t, t \in \mathbb{N}\}$  satisfies the usual consistency conditions. We conclude by applying the Kolmogorov extension theorem.  $\blacksquare$

For  $x \in X$ , we use the shorthand notation  $\mathbb{P}_x = \mathbb{P}_{\delta_x}$ . We denote by  $\mathbb{E}_\nu$  the expectation associated to  $\mathbb{P}_\nu$ . For all  $A \in \mathcal{F}$ , the function  $x \mapsto \mathbb{P}_x(A)$  is  $\mathcal{X}$ -measurable. For all  $\nu \in \mathbb{M}_1(\mathcal{X})$  and  $A \in \mathcal{F}$ ,  $\mathbb{P}_\nu(A) = \int_X \mathbb{P}_x(A) \nu(dx)$ .

### 2.4 Invariant measures

Stochastic processes are *strict-sense stationary* if, for any integer  $k$ , the distribution of the random vector  $(X_t, \dots, X_{t+k})$  does not depend on the time-shift  $t$ . In general, a Markov chain will not be stationary. For example, suppose we define an AR(1) model  $\{X_t, t \in \mathbb{N}\}$ , as  $X_t = \phi X_{t-1} + Z_t$  where  $|\phi| < 1$  and  $Z_t \sim_{\text{iid}} N(0, \sigma^2)$ . To initialize the process, set for example  $X_0 = Z_0$ . In this case,  $X_0 \sim N(0, \sigma^2)$ , but  $X_1 = \phi X_0 + Z_1 \sim N(0, \sigma^2(1 + \phi^2))$ , so that  $X_0$  and  $X_1$  do not have the same distribution. A simple fix is to put  $X_0 = Z_0 / \sqrt{1 - \phi^2}$ , in which case

$$X_t = \phi^t X_0 + \sum_{s=0}^{t-1} \phi^s Z_{t-s} \sim N(0, \sigma^2 / (1 - \phi^2)),$$

for all  $t \geq 0$ .

Under appropriate conditions on the Markov kernel  $P$ , it is possible to produce a stationary process with a proper choice of the initial distribution  $\pi$ . Assuming that such a distribution exists, the stationarity of the marginal distribution implies that  $\mathbb{E}_\pi[\mathbb{1}_A(X_0)] = \mathbb{E}_\pi[\mathbb{1}_A(X_1)]$  for any  $A \in \mathcal{X}$ . This can equivalently be written as  $\pi(A) = \pi P(A)$ , or  $\pi = \pi P$ . We can iterate to give, for any  $h \in \mathbb{N}$ ,  $\pi P^h = \pi$ . For all integers  $h$  and  $n$ , and all  $A_0, \dots, A_n \in \mathcal{X}$ ,

$$\begin{aligned} \mathbb{P}_\pi(X_h \in A_0, X_{h+1} \in A_1, \dots, X_{h+n} \in A_n) &= \int \dots \int \pi P^h(dx_0) \prod_{i=1}^n P(x_{i-1}, dx_i) \mathbb{1}_{A_i}(x_i) \\ &= \int \dots \int \pi(dx_0) \prod_{i=1}^n P(x_{i-1}, dx_i) \mathbb{1}_{A_i}(x_i) = \mathbb{P}_\pi(X_0 \in A_0, X_1 \in A_1, \dots, X_n \in A_n), \end{aligned}$$



which shows that, for any integers  $h$  and  $n$ , the random vectors  $(X_h, X_{h+1}, \dots, X_{h+n})$  and  $(X_0, X_1, \dots, X_n)$  have the same distributions. Therefore, the Markov property implies that all finite-dimensional distributions of  $\{X_t, t \in \mathbb{N}\}$  are also invariant under translation in time. These considerations lead to the definition of the *invariant measure*.

**Definition 2.15 (Invariant measure).** Given a Markov kernel  $P$ , a  $\sigma$ -finite measure  $\pi$  on  $(X, \mathcal{X})$  with the property

$$\pi(A) = \pi P(A), \quad A \in \mathcal{X},$$

will be called *invariant* (with respect to  $P$ ).

If an invariant measure is finite, it may be normalized to an *invariant probability measure*. In general, there may exist more than one invariant measure, and if  $X$  is not finite, an invariant measure may not exist. As a trivial example, consider  $X = \mathbb{N}$  and  $P(x, x+1) = 1$ .

It turns out that we can extend a stationary Markov chain to have time  $t$  take on negative values, as well.

**Proposition 2.16** *Let  $(X, \mathcal{X})$  be a measurable space. Let  $P$  a Markov kernel on  $(X, \mathcal{X})$  admitting  $\pi$  has an invariant distribution. Then, there exists on some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  a stochastic process  $\{X_t\}$  such that, for any integer  $p$ , any  $p$ -tuple  $t_1 < t_2 < \dots < t_p$ , and any  $A \in \mathcal{X}^{\otimes p}$ ,*

$$\begin{aligned} \mathbb{P}((X_{t_1}, X_{t_2}, \dots, X_{t_p}) \in A) \\ = \int \dots \int \pi(dx_1) P^{t_2-t_1}(x_1, dx_2) \dots P^{t_p-t_{p-1}}(x_{p-1}, dx_p) \mathbb{1}_A(x_1, \dots, x_p). \end{aligned} \quad (2.10)$$

PROOF. This is a direct consequence of the Kolmogorov extension theorem. ■

**Remark 2.17** *We may take  $\Omega = \prod_{t=-\infty}^{\infty} X_t$  and  $\mathcal{F} = \bigvee_{t=-\infty}^{\infty} \mathcal{X}_t$ , where  $(X_t, \mathcal{X}_t)$  is a copy of  $(X, \mathcal{X})$ . For  $t \in \mathbb{Z}$ ,  $X_t$  is the  $t$ -th coordinate mapping, i.e. for  $\{\omega\}$  an element of  $\Omega$ ,  $X_t(\omega) = \omega_t$ . This is the canonical two-sided extension of a stationary Markov chain.*

**Example 2.18 (Markov chain over a finite state-space)** *Consider a Markov chain on a finite state space  $X = \{1, \dots, n\}$  with transition kernel  $P$ . A probability measure is a vector  $\xi$  with nonnegative entries summing to 1,  $\sum_x \xi(x) = 1$ . If  $\xi$  is the initial distribution, then after one step, the distribution of the chain is  $\xi P$ ; after  $t$  steps, the distribution is  $\xi P^t$ . The probability  $\pi$  is stationary if and only if  $\pi P = \pi$ . This means that 1 should be an eigenvalue of  $P$  and in this case  $\pi$  is the left-eigenvector of  $P$  associated to the eigenvalue 1. It may be shown that, provided there exists an integer  $m$  such that  $P^m(x, x') > 0$  for all  $(x, x') \in X \times X$  (the Markov kernel  $P$  is said to be irreducible in such a case), such distribution exists and is unique.*

**Example 2.19 (Gaussian AR(1) processes)** *Consider a Gaussian AR(1) process,  $X_t = \mu + \phi X_{t-1} + \sigma Z_t$ , where  $\{Z_t, t \in \mathbb{N}\}$  is an i.i.d. sequence of standard Gaussian random variables, independent of  $X_0$ . Assume that  $|\phi| < 1$  and that  $X_0$  is Gaussian with mean  $\mu_0$  and variance  $\gamma_0^2$ . Then  $X_1$  is Gaussian with mean  $\mu + \phi \mu_0$ , and variance  $\phi^2 \gamma_0^2 + \sigma^2$ . If we choose*

$$\begin{cases} \mu + \phi \mu_0 = \mu_0 \\ \phi^2 \gamma_0^2 + \sigma^2 = \gamma_0^2 \end{cases} \quad \Rightarrow \quad \begin{cases} \mu_0 = \mu / (1 - \phi) \\ \gamma_0^2 = \sigma^2 / (1 - \phi^2) \end{cases}$$

*then  $X_1$  and  $X_0$  have the same distribution. Therefore, the Gaussian distribution with mean  $\mu / (1 - \phi)$  and variance  $\sigma^2 / (1 - \phi^2)$  is a stationary distribution. We will show later that this distribution is unique.*

**Example 2.20** Consider a Markov chain whose state space  $X = (0, 1)$  is the open unit interval. If the chain is at  $x$ , it picks one of the two intervals  $(0, x)$  or  $(x, 1)$  with equal probability  $1/2$ , and then moves to a point  $y$  which is uniformly distributed in the chosen interval. This Markov chain has a transition density with respect to Lebesgue measure on the interval  $(0, 1)$ , which is given by

$$k(x, y) = \frac{1}{2} \frac{1}{x} \mathbb{1}_{(0, x)}(y) + \frac{1}{2} \frac{1}{1-x} \mathbb{1}_{(x, 1)}(y). \quad (2.11)$$

The first term in the sum corresponds to a move from  $x$  to the interval  $(0, x)$ ; the second, to a move from  $x$  to the interval  $(x, 1)$ .

This Markov chain can be equivalently represented as an iterated random sequence. Let  $\{U_t, t \in \mathbb{N}\}$  be a sequence of i.i.d. random variable uniformly distributed on the interval  $(0, 1)$ . Let  $\{\varepsilon_t, t \in \mathbb{N}\}$  be a sequence of i.i.d. random Bernoulli variables distributed with probability of success  $1/2$ , independent of  $\{U_t, t \in \mathbb{N}\}$ . Let  $X_0$ , the initial state, be distributed according to some initial distribution  $\xi$  on  $(0, 1)$ , and be independent of  $\{U_t, t \in \mathbb{N}\}$  and  $\{\varepsilon_t, t \in \mathbb{N}\}$ . Define the sequence  $\{X_t, t \in \mathbb{N}\}$  for  $t \geq 1$  as follows

$$X_t = \varepsilon_t [X_{t-1} U_t] + (1 - \varepsilon_t) [X_{t-1} + U_t(1 - X_{t-1})], \quad (2.12)$$

Begin by assuming that the stationary distribution has a density  $p$  with respect to Lebesgue measure. From (2.11),

$$p(y) = \int_0^1 k(x, y) p(x) dx = \frac{1}{2} \int_y^1 \frac{p(x)}{x} dx + \frac{1}{2} \int_0^y \frac{p(x)}{1-x} dx. \quad (2.13)$$

Differentiation gives

$$p'(x) = -\frac{1}{2} \frac{p(x)}{x} + \frac{1}{2} \frac{p(x)}{1-x} \quad \text{or} \quad \frac{p'(x)}{p(x)} = \frac{1}{2} \left( -\frac{1}{x} + \frac{1}{1-x} \right).$$

The solutions of this first-order non-linear differential equation are given by

$$f_C(x) = \frac{C}{\sqrt{x(1-x)}}, \quad (2.14)$$

where  $C \in \mathbb{R}$  is a constant. Note that for  $x \in (0, 1)$  and  $C \geq 0$ ,  $f_C(x) \geq 0$ . For all  $z \in (0, 1)$ ,

$$\int_0^z f_C(y) dy = 2C \arcsin(\sqrt{z}). \quad (2.15)$$

Choosing  $C = 1/\pi$  makes  $\int_0^1 f_{1/\pi}(y) dy = 1$ . Therefore, the arcsine density  $p = f_{1/\pi}$  is a stationary distribution. It is the unique stationary distribution admitting a density with respect to Lebesgue measure. We will show later that  $p$  is indeed the unique stationary distribution.

## 2.5 Observation-driven models

Of course there are many time series models that are not Markov chains. However, in many useful situations, Markov chains are not very far from the surface.

**Definition 2.21 (Observation-driven model).** Let  $(X, \mathcal{X})$  and  $(Y, \mathcal{Y})$  be a measurable space,  $Q$  be Markov kernel on  $X \times \mathcal{Y}$  and  $(x, y) \mapsto f_y(x)$  a measurable function from  $(X \times Y, \mathcal{X} \otimes \mathcal{Y})$  to  $(X, \mathcal{X})$ . Let  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}, \mathbb{P})$  be a filtered probability space.

An *observation-driven time series model* is a stochastic process  $\{(X_t, Y_t), t \in \mathbb{N}\}$  adapted to  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}, \mathbb{P})$  taking values in  $X \times Y$  satisfying the following recursions: for all  $t \in \mathbb{N}^*$ ,

$$\mathbb{P}[Y_t \in A | \mathcal{F}_{t-1}] = \mathbb{P}[Y_t \in A | X_{t-1}] = Q(X_{t-1}, A), \quad \text{for any } A \in \mathcal{Y}, \quad (2.16)$$

$$X_t = f_{Y_t}(X_{t-1}). \quad (2.17)$$

The name *observation-driven* models was introduced in Cox (1981), but the definition used here is slightly different from the one used in this original contribution. For any integer  $t > 1$ ,  $X_{t-1} = f_{Y_{t-1}} \circ f_{Y_{t-2}} \circ \dots \circ f_{Y_1}(X_0)$ : therefore  $X_{t-1}$  is a function of the trajectory, up to time  $t-1$  and the initial condition  $X_0$ . The state  $X_{t-1}$  summarizes the information available on the conditional distribution of  $Y_t$  given  $X_0, Y_1, \dots, Y_{t-1}$ . For any  $C \in \mathcal{X} \otimes \mathcal{Y}$ , and  $t \geq 1$ , we get

$$\mathbb{E}[\mathbb{1}_C(X_t, Y_t) | \mathcal{F}_{t-1}] = \mathbb{E}[\mathbb{1}_C(f_{Y_t}(X_{t-1})) | \mathcal{F}_{t-1}] = \int Q(X_{t-1}, dy) \mathbb{1}_C(f_y(X_{t-1})),$$

showing that  $\{(X_t, Y_t), t \geq 0\}$  is a Markov chain on the product space  $\mathbf{X} \times \mathbf{Y}$ , with transition kernel  $P((x, y), C) = \int Q(x, dy') \mathbb{1}_C(f_{y'}(x), y')$ , for any  $C \in \mathcal{X} \otimes \mathcal{Y}$ . Note that this transition kernel depends only upon  $x \in \mathbf{X}$ . For any  $A \in \mathcal{X}$ , we may write similarly,

$$\mathbb{E}[\mathbb{1}_A(X_t) | \mathcal{F}_{t-1}] = \mathbb{E}[\mathbb{1}_A(f_{Y_t}(X_{t-1})) | \mathcal{F}_{t-1}] = \int Q(X_{t-1}, dy) \mathbb{1}_A(f_y(X_{t-1})),$$

showing that  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is also a Markov chain on  $\mathbf{X}$ , with transition kernel  $H(x, A) = \int Q(x, dy) \mathbb{1}_A(f_y(x))$ , for any  $x \in \mathbf{X}$  and  $A \in \mathcal{X}$ .

**Example 2.22 (ARMA(1,1) model)** Consider for example an ARMA(1,1) model,

$$Y_t - \phi_1 Y_{t-1} = Z_t + \theta_1 Z_{t-1}, \quad t \geq 1, \quad (2.18)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is a sequence of i.i.d. random variables with density  $q$  with respect to the Lebesgue measure on  $\mathbb{R}$ , and  $\{Z_t, t \in \mathbb{N}\}$  is independent of  $Y_0$ , which has distribution  $\xi$ . The process of interest,  $\{Y_t, t \in \mathbb{N}\}$  is referred to as the observations. Consider  $X_t = [Y_t, Z_t]' \in \mathbf{X} = \mathbb{R}^2$ . Since  $Y_t = [\phi_1, \theta_1]X_{t-1} + Z_t$ , the conditional distribution of  $Y_t$  given  $\mathcal{F}_{t-1} = \sigma((X_s, Y_s), 0 \leq s \leq t-1)$  is given for any  $A \in \mathcal{B}(\mathbb{R})$  by

$$\mathbb{P}[Y_t \in A | \mathcal{F}_{t-1}] = \int_A q(y - [\phi_1, \theta_1]X_{t-1}) \lambda(dy),$$

and

$$X_t = \begin{bmatrix} Y_t \\ Z_t \end{bmatrix} = \begin{bmatrix} Y_t \\ Y_t \end{bmatrix} - \begin{pmatrix} 0 & 0 \\ \phi_1 & \theta_1 \end{pmatrix} X_{t-1}.$$

**Example 2.23 (GARCH(1,1) model)** The GARCH(1,1) model is defined as

$$Y_t = \sigma_t \varepsilon_t, \quad \sigma_t^2 = \alpha_0 + \beta_1 \sigma_{t-1}^2 + \alpha_1 Y_{t-1}^2, \quad t \geq 1, \quad (2.19)$$

where the driving noise  $\{\varepsilon_t, t \in \mathbb{N}\}$  is an i.i.d. sequence,  $\sigma_0^2$  is a random variable distributed according to some initial distribution  $\xi$  on  $(\mathbb{R}^+, \mathcal{B}(\mathbb{R}^+))$  assumed to be independent of  $\{\varepsilon_t, t \in \mathbb{N}\}$ . The case  $\beta_1 = 0$  corresponds to an ARCH(1) model.

The GARCH(1,1) model can be cast into the framework of observation-driven models by setting  $X_t = \sigma_{t+1}^2$  and  $f_y(x) = \alpha_0 + \alpha_1 y^2 + \beta_1 x$ .

The square volatility  $\{\sigma_t^2, t \in \mathbb{N}\}$  satisfies the following stochastic recurrence equation

$$\sigma_t^2 = \alpha_0 + \beta_1 \sigma_{t-1}^2 + \alpha_1 \sigma_{t-1}^2 \varepsilon_{t-1}^2 = \alpha_0 + (\beta_1 + \alpha_1 \varepsilon_{t-1}^2) \sigma_{t-1}^2, \quad t \geq 1, \quad (2.20)$$

By construction,  $\{\sigma_t^2, t \in \mathbb{N}\}$  is a Markov chain. By iterating the relation  $\sigma_t^2 = \alpha_0 + \beta_1 \sigma_{t-1}^2 + \alpha_1 Y_{t-1}^2$  backwards in time, we get

$$\sigma_t^2 = \sum_{j=0}^{t-1} \beta_1^j (\alpha_0 + \alpha_1 Y_{t-j-1}^2) + \beta_1^t \sigma_0^2, \quad (2.21)$$

showing that the state  $\sigma_t^2$  at time  $t$  is a deterministic function of the lagged observations  $Y_0, \dots, Y_{t-1}$  and of the initial state  $\sigma_0^2$ .

## 2.6 Iterated random functions

### 2.6.1 Strict stationarity

Let  $(X, d)$  be a Polish space and denote by  $\mathcal{X}$  the associated Borel  $\sigma$ -field. Consider now  $\{X_t, t \in \mathbb{Z}\}$  a stochastic process on  $(X, \mathcal{X})$  satisfying the following recurrence equation: for all  $t \in \mathbb{Z}$ ,

$$X_t = f_{Z_t}(X_{t-1}), \quad (2.22)$$

where  $\{Z_t, t \in \mathbb{Z}\}$  is a stochastic process on  $(Z, \mathcal{Z}^{\otimes \mathbb{Z}})$  where  $(Z, \mathcal{Z})$  a measurable space. Moreover, throughout this section, the function  $(z, x) \mapsto f_z(x)$  is assumed to be  $(\mathcal{Z} \otimes \mathcal{X}, \mathcal{X})$ -measurable. We will need the following additional assumptions.

**H 2.24** *The sequence  $\{Z\}$  is strict-sense stationary and ergodic.*

**H 2.25** *There exists a measurable function  $z \mapsto K_z$  such that*

$$d(f_z(x), f_z(y)) \leq K_z d(x, y), \quad (2.23)$$

*for all  $(x, y, z) \in X \times X \times Z$  and*

$$\mathbb{E} [\ln^+(K_{Z_0})] < \infty \quad \text{and} \quad \mathbb{E} [\ln(K_{Z_0})] < 0. \quad (2.24)$$

**H 2.26** *There exists  $x_0 \in X$  such that*

$$\mathbb{E} [\ln^+ d(x_0, f_{Z_0}(x_0))] < \infty. \quad (2.25)$$

**Remark 2.27** *Assume H 2.25: for all  $(x_0, x'_0) \in X \times X$  and  $z \in Z$ , we have*

$$d(x'_0, f_z(x'_0)) \leq (1 + K_z) d(x_0, x'_0) + d(x_0, f_z(x_0)),$$

*so that*

$$1 \wedge d(x'_0, f_z(x'_0)) \leq 1 \wedge [(1 + K_z) d(x_0, x'_0)] + 1 \wedge d(x_0, f_z(x_0)).$$

*Taking the logarithm function on both sides of the inequality and using that  $\ln(x+y) \leq \ln(x) + \ln(y)$  for all  $x, y \geq 2$ , we obtain that if H 2.26 holds for some  $x_0 \in X$ , then it also holds for all  $x_0 \in X$ .*

**Theorem 2.28.** *Assume H 2.24, H 2.25 and H 2.26. Then, for all  $x_0 \in X$ ,*

$$X_{k,t}(x_0) = f_{Z_t} \circ \dots \circ f_{Z_{t-k+1}}(x_0), \quad t \in \mathbb{Z}, \quad (2.26)$$

converges  $\mathbb{P}$ -a.s. to a random variable  $\tilde{X}_t$  which does not depend on  $x_0$ . Moreover,  $\{\tilde{X}_t, t \in \mathbb{Z}\}$  is the only strict-sense stationary solution of (2.22). Furthermore, for an arbitrary random variable  $U$ ,

$$\lim_{t \rightarrow \infty} d(X_t(U), \tilde{X}_t) = 0, \quad \mathbb{P} - \text{a.s.},$$

where  $X_t(U)$  satisfies (2.22) for all  $t \in \mathbb{N}$  and  $X_0(U) = U$ . In particular, the random variable  $X_t(U)$  converges in distribution to  $\tilde{X}_0$ , i.e., as  $t \rightarrow \infty$ ,  $X_t(U) \xrightarrow{\mathcal{L}\mathbb{P}} \tilde{X}_0$ .

A technical lemma is first needed.

**Lemma 2.29** *Let  $\{(A_t, B_t), t \in \mathbb{Z}\}$  be a strict-sense stationary and ergodic real-valued stochastic process. Assume that*

$$\mathbb{E}[\ln^+ |A_0|] < \infty, \quad -\infty \leq \mathbb{E}[\ln |A_0|] < 0, \quad \text{and} \quad \mathbb{E}[\ln^+ |B_0|] < \infty, \quad (2.27)$$

where  $x^+ = \max(0, x)$  for  $x \in \mathbb{R}$ .

$$\sum_{j=0}^{\infty} \left| \left( \prod_{i=t-j}^{t-1} A_i \right) B_{t-1-j} \right| < \infty, \quad \mathbb{P} - \text{a.s.} \quad (2.28)$$

PROOF. We will show that

$$\limsup_{j \rightarrow \infty} (|A_{t-1} \cdots A_{t-j}| |B_{t-1-j}|)^{1/j} < 1, \quad \mathbb{P} - \text{a.s.} \quad (2.29)$$

The proof then follows by application of the Cauchy root criterion for the absolute convergence of series.

We now prove (2.29) for a fixed  $t \in \mathbb{Z}$ . Note that the ergodicity of the process  $\{(A, B), \in\}$  implies that

$$\limsup_{j \rightarrow \infty} \frac{1}{j} \sum_{i=1}^j \ln |A_{t-i}| < 0, \quad \mathbb{P} - \text{a.s.} \quad (2.30)$$

Furthermore, note that by stationarity, for any  $\delta > 0$ ,

$$\begin{aligned} \sum_{j=1}^{\infty} \mathbb{P}(j^{-1} \ln^+ |B_{t-1-j}| \geq \delta) &= \sum_{j=1}^{\infty} \mathbb{P}(\delta^{-1} \ln^+ |B_0| \geq j) \\ &= \mathbb{E} \left( \sum_{j=1}^{\infty} \mathbb{1}_{\{j \leq \delta^{-1} \ln^+ |B_0|\}} \right) \leq \delta^{-1} \mathbb{E}(\ln^+ |B_0|) < \infty, \end{aligned}$$

and the Borel-Cantelli Lemma therefore shows that  $\limsup_{j \rightarrow \infty} j^{-1} \ln^+ |B_{t-1-j}| = 0$ ,  $\mathbb{P} - \text{a.s.}$  Combining this with (2.30), we get

$$\limsup_{j \rightarrow \infty} \ln (|A_{t-1}| \cdots |A_{t-j}| |B_{t-1-j}|)^{1/j} < 0, \quad \mathbb{P} - \text{a.s.}$$

which implies (2.29) by exponentiating the two terms. ■

PROOF. [of Theorem 2.28] For all  $x \in \mathbf{X}$  and all  $(k, t) \in \mathbb{N}^* \times \mathbb{Z}$ , define

$$X_{k,t}(x) := f_{Z_t} \circ \cdots \circ f_{Z_{t-k+1}}(x),$$

with the convention  $X_{0,t}(x) = x$ . According to (2.23), we have for all  $x, x' \in \mathbf{X}$  and all  $(k, t) \in \mathbb{N} \times \mathbb{Z}$ ,

$$d(X_{k,t}(x), X_{k,t}(x')) \leq K_{Z_t} d(X_{k-1,t-1}(x), X_{k-1,t-1}(x')),$$

A straightforward induction yields

$$d(X_{k,t}(x), X_{k,t}(x')) \leq d(x, x') \prod_{s=t-k+1}^t K_{Z_s}. \quad (2.31)$$

Let  $x_0 \in \mathbf{X}$  such that  $\mathbb{E}[\ln^+ d(x_0, f_{Z_0}(x_0))] < \infty$ . By setting  $x = x_0$  and  $x' = f_{Z_{t-k}}(x_0)$  in (2.31), we obtain

$$d(X_{k,t}(x_0), X_{k+1,t}(x_0)) \leq d(x_0, f_{Z_{t-k}}(x_0)) \prod_{s=t-k+1}^t K_{Z_s}. \quad (2.32)$$

Assumptions **H** 2.24, **H** 2.25 and **H** 2.26 allow to apply Lemma 2.29. Thus, the series  $\sum_{k \geq 0} d(X_{k,t}(x_0), X_{k+1,t}(x_0))$  converges  $\mathbb{P}$ -a.s. and since  $(X, d)$  is complete, there exists a  $\mathbb{P}$ -a.s. finite random variable  $\tilde{X}_t(x_0)$  such that for all  $t \in \mathbb{N}$ ,

$$\lim_{k \rightarrow \infty} X_{k,t}(x_0) = \tilde{X}_t(x_0), \quad \mathbb{P} - \text{a.s.} \quad (2.33)$$

Moreover, using again (2.31) and applying again Lemma 2.29, we obtain that  $\tilde{X}_t(x_0)$  does not depend on  $x_0$  provided that  $\mathbb{E}[\ln^+ d(x_0, f_{Z_0}(x_0))] < \infty$  and thus, we can write  $\tilde{X}_t$  instead of  $\tilde{X}_t(x_0)$ . Since  $X_{k,t}(x_0) = f_{Z_t}(X_{k-1,t-1}(x_0))$ , we have

$$\begin{aligned} d(\tilde{X}_t, f_{Z_t}(\tilde{X}_{t-1})) &\leq d(\tilde{X}_t, X_{k,t}(x_0)) + d(X_{k,t}(x_0), f_{Z_t}(\tilde{X}_{t-1})) \\ &\leq d(\tilde{X}_t, X_{k,t}(x_0)) + K_{Z_t} d(X_{k-1,t-1}(x_0), \tilde{X}_{t-1}). \end{aligned}$$

Using (2.33), we obtain that the right-hand side converges  $\mathbb{P}$ -a.s. to 0 as  $k$  goes to infinity. Thus, for all  $t \in \mathbb{Z}$ ,

$$\tilde{X}_t = f_{Z_t}(\tilde{X}_{t-1}), \quad \mathbb{P} - \text{a.s.}$$

Moreover, let  $h, t_1, \dots, t_p \in \mathbb{N}$ . Then, by stationarity of the process  $\{A_t, t \in \mathbb{Z}\}$ , we have that for all  $k \in \mathbb{N}$ ,  $(X_{k,t_1}(x_0), \dots, X_{k,t_p}(x_0))$  has the same distribution as  $(X_{k,t_1+h}(x_0), \dots, X_{k,t_p+h}(x_0))$  so that by letting  $k$  go to infinity,  $(\tilde{X}_{t_1}, \dots, \tilde{X}_{t_p})$  has the same distribution as  $(\tilde{X}_{t_1+h}, \dots, \tilde{X}_{t_p+h})$ , showing that the process  $\{\tilde{X}_t, t \in \mathbb{N}\}$  is strict-sense stationary.

Let  $\{\check{X}_t, t \in \mathbb{Z}\}$  be another strict-sense stationary solution of (2.22). Then, applying again (2.31),

$$d(\tilde{X}_t, \check{X}_t) \leq d(\tilde{X}_{t-k}, \check{X}_{t-k}) \prod_{\ell=t-k+1}^t K_{Z_\ell},$$

so that for all  $M > 0$  and all  $x_0 \in X$ ,

$$\begin{aligned} \mathbb{P}(d(\tilde{X}_t, \check{X}_t) > \varepsilon) &\leq \mathbb{P}(d(\tilde{X}_{t-k}, \check{X}_{t-k}) > M) + \mathbb{P}\left(M \prod_{\ell=t-k+1}^t K_{Z_\ell} > \varepsilon\right) \\ &\leq \mathbb{P}(d(\tilde{X}_0, x_0) > M/2) + \mathbb{P}(d(x_0, \check{X}_0) > M/2) + \mathbb{P}\left(M \prod_{\ell=t-k+1}^t K_{Z_\ell} > \varepsilon\right), \end{aligned}$$

where we have used the triangle inequality and the stationarity of the processes  $\{\tilde{X}_t, t \in \mathbb{N}\}$  and  $\{\check{X}_t, t \in \mathbb{Z}\}$ . Now again, using Lemma 2.29, we let  $k$  go to infinity to obtain

$$\mathbb{P}(d(\tilde{X}_t, \check{X}_t) > \varepsilon) \leq \mathbb{P}(d(\tilde{X}_0, x_0) > M/2) + \mathbb{P}(d(x_0, \check{X}_0) > M/2).$$

Finally, since  $M$  and  $\varepsilon$  are arbitrary,  $\tilde{X}_t$  is equal to  $\check{X}_t$   $\mathbb{P}$ -a.s., so that  $\{\tilde{X}_t, t \in \mathbb{Z}\}$  is the only strict-sense stationary solution of (2.22). Now, if  $X_t(U)$  satisfies (2.22) for all  $t \in \mathbb{N}$  and  $X_0(U) = U$ , then,

$$d(X_t(U), \tilde{X}_t) \leq d(U, \tilde{X}_0) \prod_{\ell=0}^t K_{Z_\ell},$$

and thus  $d(X_t(U), \tilde{X}_t) \xrightarrow{\mathbb{P}\text{-a.s.}} 0$  as  $t$  goes to infinity. The proof is completed.  $\blacksquare$

We now generalize Theorem 2.28. The proof closely follows the lines of the proof of Theorem 2.28 and is left as an exercise for the reader.

**Theorem 2.30.** *Let  $(z_1, \dots, z_p) \mapsto K_{z_{1:p}}$  be a  $(\mathcal{X}^{\otimes p}, \mathcal{B}(\mathbb{R}))$ -measurable function from  $Z^p$  to  $\mathbb{R}$  and let  $z \mapsto L_z$  be a  $(\mathcal{X}, \mathcal{B}(\mathbb{R}))$ -measurable function from  $Z$  to  $\mathbb{R}$ . Assume that for all  $x, x' \in X$  and all  $z_{1:p} = (z_1, \dots, z_p) \in Z^p$ ,*

$$d(f_{z_p} \circ \dots \circ f_{z_1}(x), f_{z_p} \circ \dots \circ f_{z_1}(x')) \leq K_{z_{1:p}} d(x, x') \quad (2.34)$$

$$d(f_z(x), f_z(x')) \leq L_z d(x, x'). \quad (2.35)$$

*In addition, assume that  $\mathbb{E}[\ln^+ L_{Z_0}] < \infty$ ,  $\mathbb{E}[\ln^+(K_{Z_{0:p-1}})] < \infty$ ,  $\mathbb{E}[\ln(K_{Z_{0:p-1}})] < 0$ , and  $\mathbb{E}[\ln^+ d(x_0, f_{Z_0}(x_0))] < \infty$  for some  $x_0 \in X$ . Then, the conclusions of Theorem 2.28 hold.*

**Example 2.31 (Random coefficient autoregression: the vector case)** Set  $X = \mathbb{R}^d$  and define the sequence  $\{Y\}$  by the following linear recursion:

$$Y_{t+1} = A_t Y_t + B_t, \quad t \in \mathbb{Z}, \quad (2.36)$$

where  $\{(A_t, \mathbf{B}_t), t \in \mathbb{Z}\}$  is a sequence of i.i.d. random elements defined on  $(\Omega, \mathcal{F}, \mathbb{P})$  taking values in  $\mathbb{M}_d(\mathbb{R}) \times \mathbb{R}^d$ , where  $\mathbb{M}_d(\mathbb{R})$  is the set of  $d \times d$  matrices with real coefficients. Let  $\|\cdot\|$  be a matrix norm on  $\mathbb{R}^d$  associated to a vector norm on  $\mathbb{R}^d$  denoted by  $\|\cdot\|$ . Suppose that

$$\mathbb{E}[\log^+ \|A_t\|] < \infty, \quad \mathbb{E}[\log^+ \|\mathbf{B}_t\|] < \infty. \quad (2.37)$$

The existence and uniqueness of a stationary solution of (2.36) is obtained under the condition that  $\gamma$ , the top-Lyapunov exponent satisfies

$$\gamma = \inf_{n>0} \frac{1}{n} \mathbb{E}[\log \|A_1 \dots A_n\|] < 0. \quad (2.38)$$

This result may also be obtained from Theorem 2.30. Write (2.36) as  $\mathbf{Y}_{t+1} = f_{Z_t}(\mathbf{Y}_t)$  with  $Z_t = (A_t, \mathbf{B}_t)$  and  $f_{(A, \mathbf{B})}(y) = Ay + \mathbf{B}$ . Then,

$$f_{Z_p} \circ \dots \circ f_{Z_1}(y) = A_p \dots A_1 y + \sum_{j=1}^p (A_{j-1} \dots A_1) \mathbf{B}_j,$$

showing that

$$\|f_{Z_p} \circ \dots \circ f_{Z_1}(y) - f_{Z_p} \circ \dots \circ f_{Z_1}(y')\| \leq \|A_p \dots A_1\| \|y - y'\|.$$

Then, Theorem 2.30 applies on the condition that for some  $p > 0$ ,

$$\mathbb{E}[\log \|A_p \dots A_1\|] < 0,$$

which is exactly equivalent to the fact that the top-Lyapunov exponent  $\gamma$  is strictly less than 0.

### 2.6.2 Weak stationarity

We now consider conditions upon which (2.22) admits moments of order  $p \geq 1$ .

**H 2.32**  $\{Z\}$  is a sequence of i.i.d. random elements.

**H 2.33** There exists  $\alpha \in (0, 1)$ ,  $p \geq 1$  such that

$$\|d(f_{Z_0}(x), f_{Z_0}(y))\|_p \leq \alpha d(x, y), \quad \text{for all } (x, y) \in \mathbb{X} \times \mathbb{X}. \quad (2.39)$$

**H 2.34** There exists  $x_0 \in \mathbb{X}$  such that

$$\|d(x_0, f_{Z_0}(x_0))\|_p < \infty. \quad (2.40)$$

**Remark 2.35** Note that, under **H 2.33**, if (2.40) holds for one  $x_0 \in \mathbb{X}$ , then it holds for all  $x'_0 \in \mathbb{X}$ . Indeed, by Minkowski's inequality, we have

$$\begin{aligned} \|d(x'_0, f_{Z_0}(x'_0))\|_p &\leq d(x_0, x'_0) + \|d(f_{Z_0}(x_0), f_{Z_0}(x'_0))\|_p + \|d(x_0, f_{Z_0}(x_0))\|_p \\ &\leq (1 + \alpha) d(x_0, x'_0) + \|d(x_0, f_{Z_0}(x_0))\|_p < \infty. \end{aligned}$$

**Theorem 2.36.** Assume **H 2.32**, **H 2.33** and **H 2.34**. Then, for all  $x_0 \in \mathbf{X}$ ,

$$X_{k,t}(x_0) = f_{Z_t} \circ \cdots \circ f_{Z_{t-k+1}}(x_0), \quad t \in \mathbb{Z}, \quad (2.41)$$

converges to a random variable  $\tilde{X}_t$ ,  $\mathbb{P}$  – a.s., which does not depend on  $x_0$  and  $\lim_{k \rightarrow \infty} \|d(X_{k,t}(x_0), \tilde{X}_t)\|_p = 0$ . Moreover,  $\{\tilde{X}_t, t \in \mathbb{Z}\}$  is the only strict-sense stationary solution of (2.22) satisfying  $\mathbb{E}[d^p(x_0, \tilde{X}_0)] < \infty$  for all  $x_0 \in \mathbf{X}$ . Furthermore, for an arbitrary random variable  $U$  such that  $\mathbb{E}[d^p(x_0, U)] < \infty$  for all  $x_0 \in \mathbf{X}$ ,

$$\lim_{t \rightarrow \infty} d(X_t(U), \tilde{X}_t) = 0, \quad \mathbb{P} - \text{a.s.},$$

where  $X_t(U)$  satisfies (2.22) for all  $t \in \mathbb{N}$  and  $X_0(U) = U$ . In particular, the random variable  $X_t(U)$  converges in distribution to  $\tilde{X}_0$ , i.e., as  $t \rightarrow \infty$ ,  $X_t(U) \xrightarrow{\mathcal{L}_\mathbb{P}} \tilde{X}_0$ .

PROOF. For all  $x \in \mathbf{X}$  and all  $(k, t) \in \mathbb{N} \times \mathbb{Z}$ , define

$$X_{k,t}(x) := f_{Z_t} \circ \cdots \circ f_{Z_{t-k+1}}(x),$$

with the convention  $X_{0,t}(x) = x$ . According to **H 2.32** and **H 2.33**, we have for all  $x, x' \in \mathbf{X}$  and all  $(k, t) \in \mathbb{N} \times \mathbb{Z}$ ,

$$\|d(X_{k,t}(x), X_{k,t}(x'))\|_p \leq \alpha \|d(X_{k-1,t-1}(x), X_{k-1,t-1}(x'))\|_p,$$

A straightforward induction yields

$$\|d(X_{k,t}(x), X_{k,t}(x'))\|_p \leq \|d(x, x')\|_p \alpha^k. \quad (2.42)$$

By setting  $x = x_0$  and  $x' = f_{Z_{t-k}}(x_0)$  in (2.42), we obtain

$$\|d(X_{k,t}(x_0), X_{k+1,t}(x_0))\|_p \leq \|d(x_0, f_{Z_{t-k}}(x_0))\|_p \alpha^k. \quad (2.43)$$

Using **H 2.34**, the series  $\sum_{k \geq 0} \|d(X_{k,t}(x_0), X_{k+1,t}(x_0))\|_p$  is converging and thus, there exists a random variable  $\tilde{X}_t(x_0)$  in  $L^p$  such that for all  $t \in \mathbb{N}$ ,

$$\lim_{k \rightarrow \infty} \|d(X_{k,t}(x_0), \tilde{X}_t(x_0))\|_p = 0, \quad \mathbb{P} - \text{a.s.} \quad (2.44)$$

Moreover, using again (2.42), we obtain that  $\tilde{X}_t(x_0)$  does not depend on  $x_0$  and thus, we can write  $\tilde{X}_t$  instead of  $\tilde{X}_t(x_0)$ . Moreover, since

$$\mathbb{E}^{1/p} \left[ \left( \sum_{k=0}^{\infty} d(X_{k,t}(x_0), X_{k+1,t}(x_0)) \right)^p \right] \leq \sum_{k=0}^{\infty} \|d(X_{k,t}(x_0), X_{k+1,t}(x_0))\|_p < \infty,$$

the series  $\sum_{k \geq 0} d(X_{k,t}(x_0), X_{k+1,t}(x_0))$  converges  $\mathbb{P}$  – a.s. and thus, there exists a  $\mathbb{P}$  – a.s. random variable  $X_{\infty,t}(x_0)$  such that  $\lim_{k \rightarrow \infty} d(X_{k,t}(x_0), X_{\infty,t}(x_0)) = 0$ . Now, by Fatou's Lemma,

$$\mathbb{E}[d(X_{\infty,t}(x_0), \tilde{X}_t)] = \mathbb{E} \left[ \liminf_{k \rightarrow \infty} d(X_{k,t}(x_0), \tilde{X}_t)^p \right] \leq \liminf_{k \rightarrow \infty} \mathbb{E}[d(X_{k,t}(x_0), \tilde{X}_t)^p] = 0,$$

which implies that  $X_{\infty,t}(x_0) = \tilde{X}_t$ .

Since  $X_{k,t}(x_0) = f_{Z_t}(X_{k-1,t-1}(x_0))$ , we have

$$\begin{aligned} \|d(\tilde{X}_t, f_{Z_t}(\tilde{X}_{t-1}))\|_p &\leq \|d(\tilde{X}_t, X_{k,t}(x_0))\|_p + \|d(X_{k,t}(x_0), f_{Z_t}(\tilde{X}_{t-1}))\|_p \\ &\leq \|d(\tilde{X}_t, X_{k,t}(x_0))\|_p + \alpha \|d(X_{k-1,t-1}(x_0), \tilde{X}_{t-1})\|_p. \end{aligned}$$

Using (2.44), we obtain that the right-hand side converges to 0 as  $k$  goes to infinity. Thus, for all  $t \in \mathbb{Z}$ ,

$$\tilde{X}_t = f_{Z_t}(\tilde{X}_{t-1}), \quad \mathbb{P} - \text{a.s.}$$

Moreover, let  $h, t_1, \dots, t_p \in \mathbb{N}$ . Then, by stationarity of the process  $\{Z_t, t \in \mathbb{Z}\}$ , we have that for all  $k \in \mathbb{N}$ ,  $(X_{k,t_1}(x_0), \dots, X_{k,t_p}(x_0))$  has the same distribution as  $(X_{k,t_1+h}(x_0), \dots, X_{k,t_p+h}(x_0))$  so that by letting  $k$  go to infinity,  $(\tilde{X}_{t_1}, \dots, \tilde{X}_{t_p})$  has the same distribution as  $(\tilde{X}_{t_1+h}, \dots, \tilde{X}_{t_p+h})$  and the process  $\{\tilde{X}_t, t \in \mathbb{N}\}$  is thus strict-sense stationary.

Let  $\{\tilde{X}_t, t \in \mathbb{Z}\}$  be another strict-sense stationary solution of (2.22) such that  $\mathbb{E}[d^p(x_0, \tilde{X}_0)] < \infty$  for all  $x_0 \in \mathbf{X}$ . Then, applying again (2.31),



$$\begin{aligned} \|d(\tilde{X}_t, \check{X}_t)\|_p &\leq \|d(\tilde{X}_{t-k}, \check{X}_{t-k})\|_p \alpha^k \\ &\leq \alpha^k \left( \|d(\tilde{X}_{t-k}, x_0)\|_p + \|d(x_0, \check{X}_{t-k})\|_p \right) = \alpha^k \left( \|d(\tilde{X}_0, x_0)\|_p + \|d(x_0, \check{X}_0)\|_p \right), \end{aligned}$$

By letting  $k$  go to infinity, we obtain that  $\tilde{X}_t$  is  $\mathbb{P}$ -a.s. equal to  $X_t$  so that  $\{\tilde{X}_t, t \in \mathbb{Z}\}$  is the only strict-sense stationary solution of (2.22) such that  $\mathbb{E}[d^p(x_0, \check{X}_0)] < \infty$  for all  $x_0 \in \mathbf{X}$ . Now, if  $X_t(U)$  satisfies (2.22) for all  $t \in \mathbb{N}$  and  $X_0(U) = U$ , then,

$$d(X_t(U), \tilde{X}_t) \leq d(U, \tilde{X}_0) \prod_{\ell=0}^t K_{Z_\ell}.$$

so that  $d(X_t(U), \tilde{X}_t) \xrightarrow{\mathbb{P}\text{-a.s.}} 0$  as  $t$  goes to infinity. The proof is completed.  $\blacksquare$

### 2.6.3 Iterated random functions

Under weak conditions on the structure of the state space, every homogeneous Markov chain  $\{X_t, t \in \mathbb{N}\}$  may be represented as a functional autoregressive process, i.e.,  $X_{t+1} = f_{Z_{t+1}}(X_t)$  where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise and  $X_0$  is independent of  $\{Z_t, t \in \mathbb{N}\}$  and is distributed according to some initial probability  $\nu$ . For simplicity, consider the case of a real valued Markov chain  $\{X_t, t \in \mathbb{N}\}$  with initial distribution  $\nu$  and Markov kernel  $P$ . Let  $X$  be a real-valued random variable and let  $F(x) = \mathbb{P}(X \leq x)$  be the cumulative distribution function of  $X$ . Let  $F^{-1}$  be the quantile function, defined as the generalized inverse of  $F$  by

$$F^{-1}(u) = \inf\{x \in \mathbb{R} : F(x) \geq u\}. \quad (2.45)$$

The right continuity of  $F$  implies that  $u \leq F(x) \Leftrightarrow F^{-1}(u) \leq x$ . Therefore, if  $U$  is uniformly distributed on  $[0, 1]$ ,  $F^{-1}(U)$  has the same distribution as  $X$ , since  $\mathbb{P}(F^{-1}(U) \leq t) = \mathbb{P}(U \leq F(t)) = F(t) = \mathbb{P}(X \leq t)$ .

Define  $F_0(t) = \nu((-\infty, t])$  and  $g = F_0^{-1}$ . Consider the function  $F$  from  $\mathbb{R} \times \mathbb{R}$  to  $[0, 1]$  defined by  $F(x, x') = P(x, (-\infty, x'])$ . Then, for all  $x \in \mathbb{R}$ ,  $F(x, \cdot)$  is a cumulative distribution function; the associated quantile function  $f(x, \cdot)$  is

$$f(x, u) = f_u(x) = \inf\{x' \in \mathbb{R} : F(x, x') \geq u\}. \quad (2.46)$$

The function  $(x, u) \mapsto f_u(x)$  is a Borel function since  $(x, x') \mapsto F(x, x')$  is itself a Borel function; see Exercise 8.25. If  $U$  is uniformly distributed on  $[0, 1]$ , then, for all  $x \in \mathbb{R}$  and  $A \in \mathcal{B}(\mathbb{R})$ , we have

$$\mathbb{P}(f_U(x) \in A) = P(x, A).$$

Let  $\{U_t, t \in \mathbb{N}\}$  be a sequence of i.i.d. random variables, uniformly distributed on  $[0, 1]$ . Define a sequence of random variables  $\{X_t, t \in \mathbb{N}\}$  by  $X_0 = g(U_0)$  and for  $t \geq 0$ ,

$$X_{t+1} = f_{U_{t+1}}(X_t).$$

Then,  $\{X_t, t \in \mathbb{N}\}$  is a Markov chain with Markov kernel  $P$  and initial distribution  $\nu$ .

**Theorem 2.37.** Assume that  $(\mathbf{X}, \mathcal{X})$  is a measurable space and that  $\mathcal{X}$  is countably generated. Let  $P$  be a Markov kernel and  $\nu$  be a probability on  $(\mathbf{X}, \mathcal{X})$ . Let  $\{U_t, t \in \mathbb{N}\}$  be a sequence of i.i.d. random variables uniformly distributed on  $[0, 1]$ . There exists a measurable application  $g$  from  $([0, 1], \mathcal{B}([0, 1]))$  to  $(\mathbf{X}, \mathcal{X})$  and a measurable application  $f$  from  $(\mathbf{X} \times [0, 1], \mathcal{X} \otimes \mathcal{B}([0, 1]))$  to  $(\mathbf{X}, \mathcal{X})$  such that the sequence  $\{X_t, t \in \mathbb{N}\}$  defined by  $X_0 = g(U_0)$  and  $X_{t+1} = f_{U_{t+1}}(X_t)$  for  $t \geq 0$ , is a Markov chain with initial distribution  $\nu$  and transition probability  $P$  defined on  $\mathbf{X} \times \mathcal{X}$  by

$$P(x, A) = \mathbb{P}(f_U(x) \in A).$$

We consider the Markov chain  $\{X_t, t \in \mathbb{N}\}$  defined by the following recurrence equation

$$X_t = f(X_{t-1}, Z_t) = f_{Z_t}(X_{t-1}), \quad t \geq 1, \quad (2.47)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise independent of the initial condition  $X_0$ . As shown in Theorem 2.37, most Markov chains can be represented in this way, but the function  $f$  will in general not display any useful property. We assume in this section that the function  $f$  has some contraction properties. We will establish, under these conditions, the existence and uniqueness of the invariant measure.

For  $x_0 \in X$ , define the *forward iteration* starting from  $X_0^{x_0} = x_0$  by

$$X_t^{x_0} = f_{Z_t}(X_{t-1}^{x_0}) = f_{Z_t} \circ \cdots \circ f_{Z_1}(x_0);$$

this is just a rewrite of equation (2.47). Now, define the *backward iteration* as

$$Y_t^{x_0} = f_{Z_1} \circ \cdots \circ f_{Z_t}(x_0). \quad (2.48)$$

Of course,  $Y_t^{x_0}$  has the same distribution as  $X_t^{x_0}$  for each  $t \in \mathbb{N}$ . We will show that, under **H 2.25-H 2.26**, the forward process  $\{X_t, t \in \mathbb{N}\}$  has markedly different behavior from the backward process  $\{Y_t^{x_0}, t \in \mathbb{N}\}$ : the forward process moves *ergodically* in  $X$ , while the backward process converges to a limit  $\mathbb{P}$ -a.s. The distribution of this limit, which does not depend on the initial state  $x_0$ , is the unique stationary distribution.

**Lemma 2.38** *Assume that there exists a  $\mathbb{P}$ -a.s. finite random variable  $Y_\infty$  such that for all  $x_0 \in X$ , the backward process  $\{Y_t^{x_0}, t \in \mathbb{N}\}$  defined in (2.48) satisfies*

$$\lim_{t \rightarrow \infty} Y_t^{x_0} = Y_\infty, \quad \mathbb{P} - \text{a.s.} \quad (2.49)$$

*Then, the Markov chain  $\{X_t, t \in \mathbb{N}\}$  defined in (2.47) with Markov kernel  $P$  implicitly defined by (2.47), admits a unique invariant distribution  $\pi$ . In addition,  $\pi$  is the distribution of  $Y_\infty$  and for any  $\xi \in \mathbb{M}_1(\mathcal{X})$ , the sequence of probability measures  $\{\xi P^t, t \in \mathbb{N}\}$  converges weakly to  $\pi$ .*

**PROOF.** Denote by  $\pi$  the distribution of  $Y_\infty$ . Since  $\{Z_t\}$  is an i.i.d. sequence, it holds that

$$Y_{t+1}^{x_0} = f_{Z_1} \circ \cdots \circ f_{Z_{t+1}}(x_0) \stackrel{\text{law}}{=} f_{Z_0} \circ \cdots \circ f_{Z_t}(x_0) = f_{Z_0}(Y_t^{x_0}).$$

Since  $f_{Z_0}$  is continuous, passing to the limit implies that  $Y_\infty \stackrel{\text{law}}{=} f_{Z_0}(Y_\infty)$ , hence  $\pi$  is an invariant probability measure for  $P$ .

We now prove that  $\pi$  is the unique invariant probability measure of  $P$ . For any  $x \in X$ , the distribution of  $X_t^x = f_{Z_t} \circ \cdots \circ f_{Z_1}(x)$  is  $\delta_x P^t$ . Since  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise, for each  $t \in \mathbb{N}$ ,  $X_t^x = f_{Z_t} \circ \cdots \circ f_{Z_1}(x)$  has the same distribution as  $Y_t^x = f_{Z_1} \circ \cdots \circ f_{Z_t}(x)$ . Thus the sequence  $\{X_t^x, t \in \mathbb{N}\}$  converges weakly to  $\pi$ , i.e., for all  $h \in C_b(X)$  and  $x \in X$ , we get

$$\lim_{t \rightarrow \infty} \delta_x P^t h = \lim_{t \rightarrow \infty} \mathbb{E}[h(X_t^x)] = \lim_{n \rightarrow \infty} \mathbb{E}[h(Y_n^x)] = \pi(h).$$

By dominated convergence, this yields, for any probability measure  $\xi$  and  $h \in C_b(X)$ ,

$$\lim_{t \rightarrow \infty} \xi P^t h = \lim_{t \rightarrow \infty} \int_X P^t h(x) \xi(dx) = \pi(h). \quad (2.50)$$

If  $\pi'$  is an invariant distribution for  $P$ , then  $\pi' P^t = \pi'$  for all  $t \in \mathbb{N}$ , so that (2.50) with  $\xi = \pi'$  yields  $\pi'(h) = \lim_{t \rightarrow \infty} \pi' P^t(h) = \pi(h)$  for any  $h \in C_b(X)$ , showing that  $\pi' = \pi$ .  $\blacksquare$

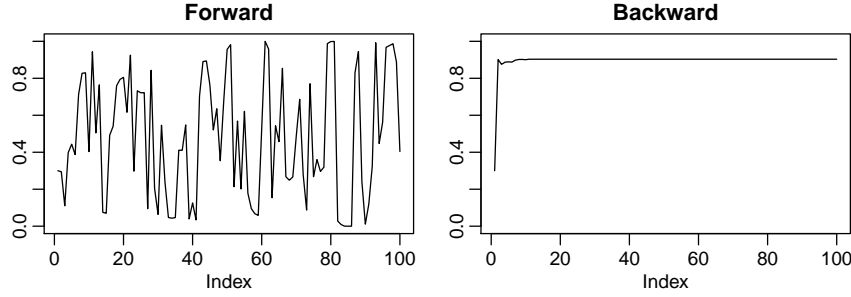
**Example 2.39** *Consider the model*

$$X_t = \varepsilon_t U_t X_{t-1} + (1 - \varepsilon_t)[X_{t-1} + U_t(1 - X_{t-1})], \quad (2.51)$$

where  $\{U_t, t \in \mathbb{N}\}$  is i.i.d. uniform on  $[0, 1]$ ,  $\{\varepsilon_t, t \in \mathbb{N}\}$  is i.i.d. Bernoulli with probability of success  $1/2$ ,  $\{U_t\}$ ,  $\{\varepsilon_t\}$  and  $X_0$  are independent. Set  $T = [0, 1] \times \{0, 1\}$  and  $\mathcal{T} = \mathcal{B}([0, 1]) \otimes \mathcal{P}\{0, 1\}$ . Then,  $X_t = f(X_{t-1}, Z_t)$  with  $Z_t = (U_t, \varepsilon_t)$  and  $f_{u,\varepsilon}(x)(u) = xu\varepsilon + (1 - \varepsilon)[x + u(1 - x)]$ . For any  $(x, y) \in [0, 1] \times [0, 1]$ ,  $|f_{u,\varepsilon}(x) - f_{u,\varepsilon}(y)| \leq K_{(u,\varepsilon)}|x - y|$  with

$$K_{u,\varepsilon} = \varepsilon u + (1 - \varepsilon)(1 - u). \quad (2.52)$$

Figure 2.1 illustrates the difference between the backward process (right hand panel, convergence) and the forward process (left hand panel, ergodic behavior). Both processes start from  $x_0 = 0.3$ , and use the same random functions to move. The order in which the functions are composed is the only difference. In the right panel, the limit 0.2257 is random because it depends on the functions being iterated; we will see in Theorem 2.40 that this limiting value does not depend on the initial condition  $x_0$ .



**Fig. 2.1** Display for Example 2.39. The left hand panel shows ergodic behavior by the forward process; the right hand side panel shows convergence of the backward process.

We now state the main result of this section, i.e., conditions upon which the Markov chain (2.47) admits a unique stationary distribution.

**Theorem 2.40.** Assume **H 2.25** and **H 2.26**. Then, the Markov chain  $\{X_t, t \in \mathbb{N}\}$  defined in (2.47) admits a unique invariant probability  $\pi$ , which is the distribution of the almost-sure limit of the sequence of random variables  $\{Y_t^{x_0}, t \in \mathbb{N}\}$ , where  $Y_t^{x_0} = f_{Z_1} \circ f_{Z_2} \circ \dots \circ f_{Z_t}(x_0)$ . In addition, for any  $\xi \in \mathbb{M}_1(\mathcal{X})$ , the sequence of probability measures  $\{\xi P^t, t \in \mathbb{N}\}$  converges weakly to  $\pi$ .

PROOF. According to Lemma 2.38, we only need to show that the backward process  $\{Y_t^{x_0}, t \in \mathbb{N}\}$  converges  $\mathbb{P}$ -a.s. to a random variable  $Y_\infty$ , which is  $\mathbb{P}$ -a.s. finite and does not depend on  $x_0$ . First note that for all  $x_0, x \in \mathcal{X}$ ,

$$\begin{aligned} d(Y_t^{x_0}, Y_t^x) &= d(f_{Z_1}[f_{Z_2} \circ \dots \circ f_{Z_t}(x_0)], f_{Z_1}[f_{Z_2} \circ \dots \circ f_{Z_t}(x)]) \\ &\leq K_{Z_1} d(f_{Z_2} \circ \dots \circ f_{Z_t}(x_0), f_{Z_2} \circ \dots \circ f_{Z_t}(x)). \end{aligned}$$

Thus, by induction,

$$d(Y_t^{x_0}, Y_t^x) \leq d(x_0, x) \prod_{i=1}^t K_{Z_i}. \quad (2.53)$$

Noting that  $Y_{t+1}^{x_0} = Y_t^x$  with  $x = f_{Z_{t+1}}(x_0)$ , we obtain

$$d(Y_t^{x_0}, Y_{t+1}^{x_0}) \leq d(x_0, f_{Z_{t+1}}(x_0)) \prod_{i=1}^t K_{Z_i}.$$

Since  $\mathbb{E}[\log K_{Z_0}] < 0$  (possibly  $\mathbb{E}[\log K_{Z_0}] = -\infty$ ), the strong law of large numbers implies that  $\limsup_{t \rightarrow \infty} t^{-1} \sum_{i=1}^t \log(K_{Z_i}) < 0$ ,  $\mathbb{P}$ -a.s. Thus,

$$\limsup_{t \rightarrow \infty} \left\{ \prod_{i=1}^t K_{Z_i} \right\}^{1/t} = \exp \left\{ \limsup_{t \rightarrow \infty} t^{-1} \sum_{i=1}^t \log(K_{Z_i}) \right\} < 1 \quad \mathbb{P}\text{-a.s.} \quad (2.54)$$

Next, since  $\mathbb{E}[\log_+(d(x_0, f_{Z_0}(x_0)))] < \infty$ , by the strong law of large numbers, it holds that  $\lim_{t \rightarrow \infty} t^{-1} \sum_{k=1}^t \log_+ d(x_0, f_{Z_k}(x_0)) < \infty$   $\mathbb{P}$ -a.s. which implies that  $\lim_{t \rightarrow \infty} t^{-1} \log_+ d(x_0, f_{Z_t}(x_0)) = 0$ ,  $\mathbb{P}$ -a.s.<sup>1</sup> Thus  $\limsup_{t \rightarrow \infty} t^{-1} \log d(x_0, f_{Z_t}(x_0)) \leq 0$   $\mathbb{P}$ -a.s. which, with (2.54), yields

<sup>1</sup> Recall that if  $\{u_k, k \in \mathbb{N}\}$  is such that  $U_n/n = \sum_{k=1}^n u_k/n$  converges as  $n$  goes to infinity, then  $u_n/n = U_n/n - [(n-1)/n]U_{n-1}/(n-1) \rightarrow 0$ .

$$\limsup_{t \rightarrow \infty} \{d(Y_t^{x_0}, Y_{t+1}^{x_0})\}^{1/t} = \limsup_{t \rightarrow \infty} \left\{ d(x_0, f(x_0, Z_{t+1})) \prod_{i=1}^t K_{Z_i} \right\}^{1/t} < 1 \quad \mathbb{P} - \text{a.s.}$$

By the Cauchy root test, this implies that the series  $\sum_{t=1}^{\infty} d(Y_t^{x_0}, Y_{t+1}^{x_0})$  is almost surely absolutely convergent. Since  $(X, d)$  is complete, this in turn implies that  $\{Y_t^{x_0}\}$  is almost surely convergent to a  $\mathbb{P} - \text{a.s.}$  finite random variable. Denote by  $Y_{\infty}^{x_0}$  the almost sure limit.

Using again (2.53) and  $\lim_{t \rightarrow \infty} \prod_{i=1}^t K_{Z_i} = 0 \quad \mathbb{P} - \text{a.s.}$ , we get that for any  $x_0, x \in X$ ,  $\lim_{t \rightarrow \infty} d(Y_t^{x_0}, Y_t^x) = 0 \quad \mathbb{P} - \text{a.s.}$  so that  $Y_{\infty}^{x_0} = Y_{\infty}^x$ . The proof then follows by applying Lemma 2.38.  $\blacksquare$

**Example 2.41 (GARCH(1, 1) model)**  $\{\sigma_t^2\}$  is a Markov chain and, for  $t \geq 1$ ,

$$\sigma_t^2 = \alpha_0 + \alpha_1 \sigma_{t-1}^2 Z_{t-1}^2 + \beta_1 \sigma_{t-1}^2 = f(\sigma_{t-1}^2, Z_{t-1}^2),$$

where  $f(x, z) = \alpha_0 + (\alpha_1 z^2 + \beta_1)x$  and  $\{Z_t, t \in \mathbb{N}\}$  is a sequence of i.i.d. random variables. Thus the GARCH(1, 1) model satisfies (2.23) with  $K_z = \alpha_1 z^2 + \beta_1$ . Consequently, a sufficient condition for the existence and uniqueness of an invariant distribution is

$$\mathbb{E} [\ln(\alpha_1 Z_1^2 + \beta_1)] < 0. \quad (2.55)$$

**Example 2.42 (Bilinear model)** Consider the bilinear process  $\{X_t, t \in \mathbb{N}\}$  defined by

$$X_t = aX_{t-1} + bZ_t X_{t-1} + Z_t, \quad (2.56)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is a sequence of integrable i.i.d. random variables and  $a$  and  $b$  are constants.

? studies this process and concludes that if  $\{Z_t\}$  is a Gaussian white noise, a sufficient condition for the existence of a stationary distribution is  $\mathbb{E}[|a + bZ_0|] < 1$ . ? relax the Gaussian assumption and simply assume that the law  $Z_0$  is absolutely continuous with respect to Lebesgue's measure on  $\mathbb{R}$ , and its density is positive and lower semicontinuous. We will show that none of these assumptions on the distribution of  $Z_0$  are necessary.

Define  $f(x, z) = (a + bz)x + z$ . The bilinear model (2.56) may be written as  $X_t = f(X_{t-1}, Z_t)$ . For any  $(x, y, z) \in \mathbb{R}$ ,  $|f(x, z) - f(y, z)| \leq |a + bz||x - y|$ , which implies that **H** 2.25 is satisfied with  $K_z = |a + bz|$  as soon as  $\mathbb{E}[\ln(|a + bZ_0|)] < 0$ . If in addition,  $\mathbb{E}[\ln_+(|Z_0|)] < \infty$ , then **H** 2.26 is also satisfied. Theorem 2.40 shows that the bilinear model (2.56) has a unique stationary distribution  $\pi$ , and that, starting from any initial distribution, the distribution of the iterates of the chain converge to  $\pi$ .

**Example 2.43 (?? cont.)** Using (2.52), we get

$$\mathbb{E}[\log(K_Z)] = (1/2)\mathbb{E}[\log(U_0)] + (1/2)\mathbb{E}[\log(1 - U_0)] = \mathbb{E}[\log(U_0)] = -1.$$

Therefore, **H** 2.25 is satisfied and the Markov chain  $\{X_t, t \in \mathbb{N}\}$  defined by (2.51) has a unique stationary distribution.

Instead of assuming that the contraction condition is satisfied pathwise, we may alternatively try to consider a contraction condition in the  $p$ -th norm.

**Theorem 2.44.** Assume **H** 2.33 and **H** 2.34. Then, the Markov chain  $\{X_t, t \in \mathbb{N}\}$  defined in (2.47) admits a unique invariant probability  $\pi$ , which is the distribution of the almost-sure limit of the sequence of random variables  $Y_t^{x_0} = f_{Z_1} \circ f_{Z_2} \circ \dots \circ f_{Z_t}(x_0)$ . In addition,

$$\int_X d^p(x_0, x) \pi(dx) < \infty.$$

For any  $\xi \in \mathbb{M}_1(\mathcal{X})$ , the sequence of probability measures  $\{\xi P^n, n \in \mathbb{N}\}$  converges weakly to  $\pi$ .

PROOF. Using Lemma 2.38, we will show that the backward process  $\{Y_t^{x_0}, t \in \mathbb{N}\}$  converges  $\mathbb{P}$ -a.s. to a random variable  $Y_\infty$  which is  $\mathbb{P}$ -a.s. finite and does not depend on  $x_0$ . Since  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise, we have

$$\begin{aligned} (Y_t^{x_0}, Y_{t+1}^{x_0}) &= (f_{Z_1}[f_{Z_2} \circ \cdots \circ f_{Z_t}(x_0)], f_{Z_1}[f_{Z_2} \circ \cdots \circ f_{Z_{t+1}}(x_0)]) \\ &\stackrel{\text{law}}{=} (f_{Z_0}[f_{Z_1} \circ \cdots \circ f_{Z_{t-1}}(x_0)], f_{Z_0}[f_{Z_1} \circ \cdots \circ f_{Z_t}(x_0)]) \\ &= (f_{Z_0}(Y_{t-1}^{x_0}), f_{Z_0}(Y_t^{x_0})). \end{aligned}$$

Applying (2.39), we obtain

$$\|d(Y_{t+1}^{x_0}, Y_t^{x_0})\|_p = \|d(f_{Z_0}(Y_{t-1}^{x_0}), f_{Z_0}(Y_t^{x_0}))\|_p \leq \alpha \|d(Y_t^{x_0}, Y_{t-1}^{x_0})\|_p.$$

Thus, by induction,  $\|d(Y_{t+1}^{x_0}, Y_t^{x_0})\|_p \leq \alpha^t \|d(x_0, f_{Z_0}(x_0))\|_p$ . By the Markov inequality, for all  $\varepsilon > 0$ ,

$$\mathbb{P}(\beta^t d(Y_{t+1}^{x_0}, Y_t^{x_0}) > \varepsilon) \leq \frac{\beta^{tp}}{\varepsilon^p} (\|d(Y_{t+1}^{x_0}, Y_t^{x_0})\|_p)^p \leq \frac{(\alpha\beta)^{tp}}{\varepsilon^p} (\|d(x_0, f_{Z_0}(x_0))\|_p)^p,$$

where  $\beta$  is chosen in  $(1, 1/\alpha)$ . The Borel-Cantelli Lemma then implies that  $\lim_{t \rightarrow \infty} \beta^t d(Y_{t+1}^{x_0}, Y_t^{x_0}) = 0$ ,  $\mathbb{P}$ -a.s. and since  $\beta > 1$ , this in its turn yields that the sequence  $\{Y_t^{x_0}, t \in \mathbb{N}\}$  is  $\mathbb{P}$ -a.s. a Cauchy sequence in  $(X, d)$ . Since  $(X, d)$  is complete, the sequence  $\{Y_t^{x_0}, t \in \mathbb{N}\}$  converges  $\mathbb{P}$ -a.s. to a random variable  $Y_\infty^{x_0}$ , which is  $\mathbb{P}$ -a.s. finite. We will now prove that  $Y_\infty^{x_0}$  does not depend on  $x_0$ . For  $x \in X$  and  $t \in \mathbb{N}$ , define  $Y_t^x = f_{Z_1} \circ \cdots \circ f_{Z_t}(x)$ . Again, applying (2.39), we obtain

$$\|d(Y_t^{x_0}, Y_t^x)\|_p = \|d(f_{Z_0}(Y_{t-1}^{x_0}), f_{Z_0}(Y_{t-1}^x))\|_p \leq \alpha \|d(Y_{t-1}^{x_0}, Y_{t-1}^x)\|_p,$$

which implies  $\|d(Y_t^{x_0}, Y_t^x)\|_p \leq \alpha^t \|d(x_0, x)\|_p$ . As above, this implies that  $d(Y_\infty^{x_0}, Y_\infty^x) = 0$ ,  $\mathbb{P}$ -a.s. Thus,  $Y_\infty^{x_0}$  does not depend on  $x_0$  and we can thus set  $Y_\infty = Y_\infty^{x_0}$ . Moreover, using  $Y_\infty = \lim_{t \rightarrow \infty} Y_t^{x_0}$ ,  $\mathbb{P}$ -a.s., we obtain, by Fatou's lemma,

$$\begin{aligned} \|d(x_0, Y_\infty)\|_p &= \left\| \liminf_{t \rightarrow \infty} d(x_0, Y_t^{x_0}) \right\|_p \leq \liminf_{t \rightarrow \infty} \|d(x_0, Y_t^{x_0})\|_p \\ &\leq \liminf_{t \rightarrow \infty} \sum_{n=0}^{t-1} \|d(Y_n^{x_0}, Y_{n+1}^{x_0})\|_p \leq (1-\alpha)^{-1} \|d(x_0, f_{Z_0}(x_0))\|_p < \infty. \end{aligned}$$

Thus,

$$\int_X d^p(x_0, x) \pi(dx) = \left( \|d(x_0, Y_\infty)\|_p \right)^p < \infty.$$

■

**Example 2.45 (GARCH(1, 1); cont.)** By Theorem 2.44, a sufficient condition for  $\mathbb{E}_\pi[\sigma_1^2] < \infty$  is  $\mathbb{E}[\alpha_1 Z_1^2 + \beta_1] < 1$ . Since we have assumed that  $\mathbb{E}[Z_1^2] = 1$ , this condition boils down to  $\alpha_1 + \beta_1 < 1$ . It is easily seen that this is also a necessary condition, since the condition  $\mathbb{E}_\pi[\sigma_1^2] < \infty$  implies that

$$\mathbb{E}_\pi[\sigma_t^2] = \alpha_0 + (\alpha_1 + \beta_1) \mathbb{E}_\pi[\sigma_{t-1}^2], \quad \mathbb{E}_\pi[\sigma_1^2] = \mathbb{E}_\pi[\sigma_t^2] = \mathbb{E}_\pi[\sigma_{t-1}^2],$$

which admits a finite solution only if  $\alpha_1 + \beta_1 < \infty$ . If, moreover,

$$\mathbb{E}_\pi[(\alpha_1 Z_1^2 + \beta_1)^p] < 1, \tag{2.57}$$

for some  $p \geq 1$ , then the stationary distribution of  $\sigma_t^2$  has a finite moment of order  $p$ . For  $p = 2$ , this condition becomes  $\alpha_1^2 \mathbb{E}[Z_1^4] + 2\alpha_1\beta_1 + \beta_1^2 < 1$ .

**Example 2.46 (Bilinear model; cont.)** Assume that there exists  $p \geq 1$  such that  $\mathbb{E}[Z_0^p] < \infty$ , then,

$$\|f(x, Z_0) - f(y, Z_0)\|_p = \|(a + bZ_0)\|_p |x - y|.$$

Thus, Assumption (2.39) holds if  $\mathbb{E}[|a + bZ_0|^p] < 1$ . In that case, (2.40) also holds and thus there exists an unique invariant probability measure  $\pi$  such that  $\mathbb{E}[|X_0|^p] < \infty$  if the distribution of  $X_0$  is  $\pi$ .

**Example 2.47 (?? cont.)** Using (2.52), for any  $p \geq 1$ , we have  $\mathbb{E}[K_{U, \varepsilon}^p] \leq 2\mathbb{E}[U^p] = 2/(p+1)$ , therefore, (2.39) is satisfied for any  $p \geq 1$ . (2.44) is satisfied showing that Theorem 2.44 is satisfied for any  $p \geq 1$ . We have already shown that the probability density  $p(x) = \sqrt{x(1-x)}/\pi$  is a stationary distribution;

Theorem 2.44 shows that this distribution is unique. In ??, we have sampled  $10^4$  independent trajectories of the process started from  $X_0 = 0.01$ . We have displayed the marginal stationary distribution after 3, 5, 10 and 20 iterations, together with the limiting distribution. In this case, the convergence to the stationary distribution is very quick.

**Example 2.48 (Log-linear Poisson autoregression)** Consider the following log-linear autoregressive models introduced by ?, defined as follows:

$$\mathcal{L}(Y_t | \mathcal{F}_{t-1}^Y) = \text{Poisson}(\exp(U_t)) \quad (2.58a)$$

$$U_t = a + bU_{t-1} + c \ln(1 + Y_{t-1}), \quad t \geq 1, \quad (2.58b)$$

where  $\mathcal{F}_t^Y = \sigma(U_0, Y_s, s \leq t)$  and  $d, a, b$  are real-valued parameters. Let  $\{N_k, k \geq 1\}$  be a sequence of independent unit rate homogeneous Poisson process on the real line, independent of  $U_0$ . Then  $\{U_n, n \in \mathbb{N}\}$  can be expressed as  $U_k = F(U_{k-1}, N_k)$ , where  $F$  is the function defined on  $\mathbb{R} \times \mathbb{N}^{\mathbb{R}}$  by

$$F(u, N) = a + bu + c \ln\{1 + N(e^u)\}.$$

The transition kernel  $P$  of the Markov chain  $\{U_t, t \in \mathbb{N}\}$  can be expressed as

$$Pf(v) = \mathbb{E}[f(a + bv + c \ln\{1 + N(e^v)\})],$$

where  $N$  is a unit rate homogeneous Poisson process.

Given the coefficients  $(a, b, c)$  and the initial intensity  $U_0$ , the log-intensity  $U_t$  can be expressed explicitly from the lagged responses by expanding (2.58b):

$$U_t = a \frac{1 - b^t}{1 - b} + b^t U_0 + c \sum_{i=0}^{t-1} b^i \ln(1 + Y_{t-i-1}).$$

Hence this model belongs to the class of observation-driven models. See ?. In this log-linear model, the lagged observation  $Y_t$  is fed into the autoregressive equation for  $U_t$  via the term  $\ln(Y_{t-1} + 1)$ . Adding one to the integer valued observation is a standard way to avoid potential problems with zero counts. In addition, both the intensity and the counts  $Y_t$  are transformed onto the same logarithmic scale. Covariates can be included in the right-hand side of (2.58b).

Let us check that **H** 2.33 holds for this model. By Minkowski's inequality, we have, for  $v \geq u \geq 0$ ,

$$\|F(u, N) - F(v, N)\|_1 \leq |b||u - v| + |c| \left\| \ln \left( \frac{1 + N(e^v)}{1 + N(e^u)} \right) \right\|_1.$$

We will prove below that for  $v \geq u \geq 0$ ,

$$\mathbb{E} \left[ \ln \left( \frac{1 + N(e^v)}{1 + N(e^u)} \right) \right] \leq v - u. \quad (2.59)$$

This yields, for  $u, v \geq 0$ ,

$$\|F(v, N) - F(u, N)\|_1 \leq (|b| + |c|)|v - u|.$$

The contraction property **H** 2.33 holds if  $|b| + |c| < 1$ .

We now prove (2.59). Since  $N$  has independent increments, we can write  $(1 + N(e^v))/(1 + N(e^u)) = 1 + W/(1 + U)$ , where  $W$  and  $U$  are independent Poisson random variable with respective means  $e^v - e^u$  and  $e^u$ . The function  $x \mapsto \ln(1 + x)$  is concave, thus, by Jensen's inequality, we obtain

$$\begin{aligned}
\mathbb{E} \left[ \ln \left( \frac{1 + N(\mathbf{e}^v)}{1 + N(\mathbf{e}^u)} \right) \right] &= \mathbb{E} \left[ \ln \left( 1 + \frac{W}{1 + U} \right) \right] \\
&\leq \ln \left( 1 + \mathbb{E} \left[ \frac{W}{1 + U} \right] \right) \\
&= \ln \left( 1 + (\mathbf{e}^v - \mathbf{e}^u) \mathbb{E} \left[ \frac{1}{1 + U} \right] \right) .
\end{aligned}$$

*Note now that*

$$\mathbb{E} \left[ \frac{1}{1 + U} \right] = \mathbf{e}^{-\mathbf{e}^u} \sum_{k=0}^{+\infty} \frac{1}{1+k} \frac{\mathbf{e}^{ku}}{k!} = \mathbf{e}^{-u} \mathbf{e}^{-\mathbf{e}^u} \sum_{k=1}^{\infty} \frac{\mathbf{e}^{ku}}{k!} \leq \mathbf{e}^{-u} .$$

*This yields*

$$\mathbb{E} \left[ \ln \left( \frac{1 + N(\mathbf{e}^v)}{1 + N(\mathbf{e}^u)} \right) \right] \leq \ln \left( 1 + (\mathbf{e}^v - \mathbf{e}^u) \mathbf{e}^{-u} \right) = v - u .$$





# Inference for Markovian Models

## 3.1 Likelihood inference

Assume that  $(X_1, \dots, X_n)$  is an observation from a collection of distributions  $(\mathbb{P}_\theta, \theta \in \Theta)$  that depends on a parameter  $\theta$  ranging over a set  $\Theta$ . A popular method for finding an estimator  $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$  is to maximize a criterion of the type  $\theta \mapsto M_n(\theta)$  over the parameter set  $\Theta$ . For notational simplicity, the dependence of  $M_n$  in the observations is implicit. Such an estimator is often called an *M-estimator*. When  $(X_1, \dots, X_n)$  are i.i.d., the criterion  $\theta \mapsto M_n(\theta)$  is often chosen to be the sample average of some known functions  $m^\theta : X \rightarrow \mathbb{R}$ ,

$$M_n(\theta) = n^{-1} \sum_{i=1}^n m^\theta(X_i). \quad (3.1)$$

Sometimes, the maximizing value is computed by setting a derivative (or the set of partial derivatives in the multidimensional case) equal to zero. In such a case, the estimator  $\hat{\theta}_n$  is defined as the solution of a system of equations of the type  $\Psi_n(\theta) = 0$ . For instance, if  $\theta$  is  $d$ -dimensional, then  $\Psi_n$  typically has  $d$  coordinate functions  $\Psi_n(\theta) = (\Psi_n^{(1)}(\theta), \dots, \Psi_n^{(d)}(\theta))$ . Such estimators are often referred to as *Z-estimators*. In the i.i.d. case,  $\Psi_n(\theta)$  is often chosen to be  $\Psi_n(\theta) = n^{-1} \sum_{i=1}^n \psi^\theta(X_i)$  where  $\psi^\theta := (\psi_1^\theta, \dots, \psi_d^\theta)$ . In such a case,  $\Psi_n(\theta) = 0$  is shorthand for the system of equations

$$n^{-1} \sum_{i=1}^n \psi_j^\theta(X_i) = 0, \quad j = 1, 2, \dots, d. \quad (3.2)$$

In many examples  $\psi_j^\theta$  is taken to be the  $j$ -th partial derivative of the function  $m_\theta$ ,  $\psi_j^\theta = \partial m^\theta / \partial \theta_j$ .

In this section, we consider *maximum likelihood estimators*. Suppose first that  $X_1, \dots, X_n$  are i.i.d. with a common density  $x \mapsto p^\theta(x)$ . Then the maximum likelihood estimator maximizes the likelihood or, equivalently, the log-likelihood, given by

$$\theta \mapsto n^{-1} \sum_{i=1}^n \ln p^\theta(X_i).$$

Thus, a maximum likelihood estimator is an M-estimator with  $m^\theta = \ln p^\theta$ . If the density is partially differentiable with respect to  $\theta$  for each  $x \in X$ , then the maximum likelihood estimator also solves (3.2), with

$$\psi^\theta(x) = \nabla \ln p^\theta(x) = \left( \frac{\partial \ln p^\theta(x)}{\partial \theta_1}, \dots, \frac{\partial \ln p^\theta(x)}{\partial \theta_d} \right)'.$$

This approach extends directly to the Markov chain context. Let  $(X, d)$  be a Polish space equipped with its Borel sigma-field  $\mathcal{X}$  and  $p$  be a positive integer. Consider  $\{Q^\theta, \theta \in \Theta\}$ , a family of Markov kernels on  $X^p \times \mathcal{X}$  indexed by  $\theta \in \Theta$  where  $(\Theta, d)$  is a compact metric space. Assume that all  $(\theta, x) \in \Theta \times X^p$ ,  $Q^\theta(x; \cdot)$

is dominated by some  $\sigma$ -finite measure  $\mu$  on  $(X, \mathcal{X})$  and denote by  $q^\theta(x; \cdot)$  its Radon-Nikodym derivative:  $q^\theta(x; y) = dQ^\theta(x; \cdot)/d\mu(y)$ . If  $\{X_t, t \in \mathbb{N}\}$  is a Markov chain of order  $p$  associated to the Markov kernel  $Q^\theta$ , then, the conditional distribution of the observations  $(X_p, \dots, X_n)$  given  $X_0, \dots, X_{p-1}$  has a density with respect to the product measure  $\mu^{\otimes(n-p+1)}$  given by

$$x_{0:n} \mapsto p^\theta(x_{p:n}|x_{0:p-1}) = \prod_{t=p}^n q^\theta(x_t - p, \dots, x_{t-1}; x_t). \quad (3.3)$$

The conditional Maximum Likelihood Estimator  $\hat{\theta}_n$ , based on the observations  $X_{0:n}$  and on the family of likelihood functions (3.3) is defined by

$$\hat{\theta}_n \in \arg \max_{\theta \in \Theta} \ln p^\theta(X_{p:n}|X_{0:p-1}). \quad (3.4)$$

Starting at  $\theta_0$ , optimization algorithms generate a sequence  $\{\theta_k, k \in \mathbb{N}\}$  that terminate when either no progress can be made or when it is apparent that a solution point has been approximated with a sufficient accuracy. Prior knowledge of the model and the data set may be used to choose  $\theta_0$  to be a reasonable estimate of the solution; otherwise the starting point must be set by the algorithm, either by a systematic approach, or in some arbitrary manner. In deciding how to move from one iterate  $\theta_k$  to the other, the algorithm will use information about the likelihood, and possibly information gathered by earlier iterates  $\theta_0, \theta_1, \dots, \theta_{k-1}$ .

A classical strategy is to use a line search: at each iteration, the algorithm chooses a direction  $\delta_k$ , and searches along this direction from the current iterate  $\theta_k$  for a new iterate with a higher function value. The distance to move along the search direction  $\delta_k$  can be found by solving approximately the following one-dimensional optimization problem

$$\max_{\alpha > 0} \ln p^{\theta_k + \alpha \delta_k}(X_{p:n}|X_{0:p-1}).$$

By solving exactly this optimization problem, we would derive the maximum benefit from moving along the direction  $\delta_k$ , but an exact maximization may be quite expensive and is usually worthless. Instead, the line search algorithm will typically generate a limited number of trial step lengths, until it finds one that loosely approximates the maximum.

The most obvious choice for search direction is the *steepest ascent* function  $\theta \mapsto \nabla \ln p^\theta(X_{p:n}|X_{0:p-1})$ :

$$\delta_k = \nabla \ln p^{\theta_k}(X_{p:n}|X_{0:p-1}).$$

The gradient  $\nabla \ln p^{\theta_k}(X_{p:n}|X_{0:p-1})$  is the *score* vector computed at  $\theta = \theta_k$ . Along all the directions we could move from  $\theta_k$ , it is the one along which  $\theta \mapsto \ln p^\theta(X_{p:n}|X_{0:p-1})$  increases most rapidly. One of the advantage of the steepest ascent algorithm is that it requires only the calculation of the score function; however, it is known to be slow on difficult problems. Line search methods may use other search directions than the gradient. In general, any direction that makes a positive angle with  $\nabla \ln p^{\theta_k}(X_{p:n}|X_{0:p-1})$  is guaranteed to produce an increase in the likelihood, provided that the step-size is chosen appropriately. Another important search direction (and perhaps the most important one of all), is the *Newton direction*. This direction is derived from the second order Taylor series approximation of log-likelihood  $\ln p^\theta(X_{p:n}|X_{0:p-1})$  in the neighborhood of the current fit  $\theta_k$ . Assuming that  $\nabla^2 \ln p^{\theta_k}(X_{p:n}|X_{0:p-1})$  is positive definite, the Newton direction is

$$\delta_k^N = - \left[ \nabla^2 \ln p^{\theta_k}(X_{p:n}|X_{0:p-1}) \right]^{-1} \nabla \ln p^{\theta_k}(X_{p:n}|X_{0:p-1}). \quad (3.5)$$

Unlike the steepest ascent algorithm, there is a “natural” step length of 1 associated with the Newton direction. Most line search implementations of Newton’s method use the unit step  $\alpha = 1$  where it is possible to adjust  $\alpha$  only when it does not produce a satisfactory increase in the value of the likelihood.

When  $\nabla^2 \ln p^{\theta_k}(X_{p:n}|X_{0:p-1})$  is not positive definite, the Newton direction may not even be defined, since  $(\nabla^2 \ln p^{\theta_k}(X_{p:n}|X_{0:p-1}))^{-1}$  may not exist. Even when it is defined, it may not satisfy the descent property  $\nabla \ln p^{\theta_k}(X_{p:n}|X_{0:p-1})' \delta_k^N < 0$ , in which case it is unsuitable as a search direction. In these situations, line

search methods modify the definition of the search direction  $\delta_k^N$  to make it satisfy the descent condition while retaining the benefit of the second-order information.

Methods that use the Newton direction have a fast rate of local convergence, typically quadratic (at least when the score and the Hessian can be estimated exactly). After a neighborhood of the solution is reached, convergence to a high accuracy solution occurs most often in a few iterations.

The main drawback of the Newton direction is the need for the Hessian (also called the *information matrix*). Explicit computation of this matrix of second derivatives can sometimes be cumbersome. Finite-difference and automatic differentiation techniques may be useful in avoiding the need to calculate second derivatives. *Quasi-Newton* search directions provide an attractive alternative to Newton's method in that they do not require computation of the Hessian and yet still attain a superlinear rate of convergence. In place of the true Hessian, Quasi-Newton strategies use an approximation that is updated after each step. This update uses the changes in the gradient to gain information about the second derivatives along the search direction.

Before going further, consider some examples:

**Example 3.1 (Discrete-valued Markov chain)** *In this case, the set  $\mathcal{X}$  is at most countable; the dominating measure  $\nu$  on  $\mathcal{X}$  is chosen to be the counting measure. If  $\mathcal{X} = \{1, \dots, K\}$  is finite and if the parameters  $\theta = (\theta_{i,j})_{K \times K}$  are the transition probabilities,  $\theta_{i,j} = q^\theta(i, j)$  for all  $(i, j) \in \{1, \dots, K\}^2$ , then the maximum likelihood estimator is given by*

$$\hat{\theta}_{n,i,j} = \frac{\sum_{t=1}^{n-1} \mathbb{1}_{i,j}(X_t, X_{t+1})}{\sum_{t=1}^n \mathbb{1}_i(X_t)}.$$

**Example 3.2 (AR models)** *Consider the AR( $p$ ) model,  $X_t = \sum_{i=1}^p \phi_i X_{t-i} + \sigma Z_t$ , where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white Gaussian noise. Here  $\theta = (\phi_1, \dots, \phi_p, \sigma^2)$  and  $\Theta$  is a compact subset of  $\mathbb{R}^p \times \mathbb{R}_+$ . The conditional log-likelihood of the observations may be written as*

$$\ln p^\theta(X_{p:n} | X_{0:p-1}) = -\frac{n-p+1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=p}^n (X_t - \sum_{j=1}^p \phi_j X_{t-j})^2.$$

*The conditional likelihood function is therefore a quadratic function in the regression parameter and a convex function in the innovation variance. The maximum likelihood estimator can be computed in such a case explicitly as follows:*

$$\begin{pmatrix} \hat{\phi}_{n,1} \\ \hat{\phi}_{n,2} \\ \vdots \\ \hat{\phi}_{n,p} \end{pmatrix} = \hat{\Gamma}_n^{-1} \begin{pmatrix} n^{-1} \sum_{t=p}^n X_t X_{t-1} \\ n^{-1} \sum_{t=p}^n X_t X_{t-2} \\ \vdots \\ n^{-1} \sum_{t=p}^n X_t X_{t-p} \end{pmatrix} \quad (3.6)$$

*where  $\hat{\Gamma}_n$  is the  $(p \times p)$  empirical covariance matrix for which the  $i, j$ -th element is defined by  $\hat{\Gamma}_n(i, j) = n^{-1} \sum_{t=p}^n X_{t-i} X_{t-j}$ . The maximum likelihood estimator for the innovation variance is given by*

$$\hat{\sigma}_n^2 = \frac{1}{n-p+1} \sum_{t=p}^n \left( X_t - \sum_{j=1}^p \hat{\phi}_{n,j} X_{t-j} \right)^2. \quad (3.7)$$

**Example 3.3 (Threshold autoregressive models)** *Consider a two-regime threshold autoregressive TAR model*

$$X_t = \left\{ \phi_{1,0} + \sum_{j=1}^{p_1} \phi_{1,j} X_{t-j} + \sigma_1 Z_t \right\} \mathbb{1}_{\{X_{t-d} \leq r\}} + \left\{ \phi_{2,0} + \sum_{j=1}^{p_2} \phi_{2,j} X_{t-j} + \sigma_2 Z_t \right\} \mathbb{1}_{\{X_{t-d} > r\}}, \quad (3.8)$$

*where  $\{Z_t, t \in \mathbb{N}\}$  is a strong Gaussian white noise with zero mean and variance 1. While the delay  $d$  may be theoretically larger than the maximum autoregressive order  $p$ , this is seldom the case in practice; hence we assume that  $d \leq p$ . The normal error assumption implies that for every  $t \geq 0$ , the conditional distribution of  $X_t$  given  $X_{t-1}, \dots, X_{t-p}$  is normal, where  $p = \max(p_1, p_2)$ . To estimate the parameters, we use the likelihood conditional to the  $p$  initial values. Assume first that the threshold parameter  $r$  and the delay parameter*

$d$  are known. In such a case, the observations may be split into two parts according to whether or not  $X_{t-d} \leq r$ . Denote by  $n_1(r, d)$  the number of observations in the lower regimes. With the observations in the lower regime, we can regress  $X_t$  on  $X_{t-1}, \dots, X_{t-p}$ , to find estimates of the autoregressive coefficients  $\{\hat{\phi}_{1,j}(r, d)\}_{j=1}^{p_1}$  and the noise variance estimate (see (3.6) and (3.7))

$$\sigma_1^2(r, d) = \frac{1}{n_1(r, d)} \sum_{t=p+1}^n \left( X_t - \sum_{j=1}^{p_1} \hat{\phi}_{1,j}(r, d) X_{t-j} \right)^2 \mathbb{I}\{X_{t-d} \leq r\}. \quad (3.9)$$

Similarly, using the  $n_2(r, d)$  observations in the upper regime (note that  $n - p = n_1(r, d) + n_2(r, d)$ ), we can find estimates of the autoregressive coefficients  $\{\hat{\phi}_{2,j}(r, d)\}_{j=1}^{p_2}$  and the noise variance,  $\sigma_2^2(r, d)$ . To estimate  $(r, d)$ , we consider the profile likelihood

$$\ell(r, d) = -\frac{n-p+1}{2} \{1 + \ln(2\pi)\} - \sum_{i=1}^2 \frac{n_i(r, d)}{2} \ln(\hat{\sigma}_i^2(r, d)), \quad (3.10)$$

which is maximized over  $r$  and  $d$ . The optimization need only be performed with  $r$  over the observations  $X_{p+1}, \dots, X_n$  and  $d \in \{1, \dots, p\}$ . Note indeed that, for a given  $d$ , the functions are constant if  $r$  lies in between two consecutive observations. To avoid inconsistent estimators, we typically restrict the search of the threshold to be between two predetermined quantiles of  $X_{p+1}, \dots, X_t$ .

**Example 3.4 (Functional autoregressive models)** Consider a functional autoregressive model:

$$X_t = a^\theta(X_{t-1}, \dots, X_{t-p}) + b^\theta(X_{t-1}, \dots, X_{t-q})Z_t, \quad (3.11)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is a strong Gaussian noise with zero mean and unit variance,  $\theta \in \Theta \subset \mathbb{R}^d$  and for each  $\theta \in \Theta$ . This model includes many of the popular non linear autoregressive models such as

(a) ARCH( $p$ )  $a^\theta(x_1, \dots, x_p) = 0$  and

$$b^\theta(x_1, \dots, x_q) = \sqrt{\alpha_0 + \alpha_1 x_1^2 + \dots + \alpha_q x_q^2}, \quad (3.12)$$

with  $\theta = (\alpha_0, \alpha_1, \dots, \alpha_q) \in \Theta$  a compact subset of  $\mathbb{R}_+^* \times \mathbb{R}_+^q$ .

(b) Autoregressive models with ARCH errors  $a^\theta(x_1, \dots, x_p) = \phi_1 x_1 + \dots + \phi_p x_p$  and

$$b^\theta(x_1, \dots, x_q) = \sqrt{\alpha_0 + \sum_{j=1}^{\ell} \alpha_j (x_j - a^\theta(x_{j+1}, x_{j+2}, \dots, x_{j+p}))^2}, \quad (3.13)$$

where  $\theta = (\phi_1, \phi_2, \dots, \phi_p, \alpha_0, \alpha_1, \dots, \alpha_\ell) \in \mathbb{R}^p \times \mathbb{R}_+^{\ell+1}$ ,  $\alpha_0 > 0$  and  $q = p + \ell$ .

(c) Logistic Smooth Transition AR( $p$ )  $b^\theta(x_1, \dots, x_p) = 1$  and

$$a^\theta(x_1, \dots, x_p) = \mu_1 + \sum_{j=1}^p \phi_{1,j} x_j + \left( \mu_2 + \sum_{j=1}^p \phi_{2,j} x_j \right) (1 + \exp(-\gamma(x_d - c)))^{-1} \quad (3.14)$$

where  $d \in \{1, \dots, j\}$  and  $\theta = (\mu_i, \phi_{i,j}, i \in \{1, 2\}, j = \{1, \dots, p\}, \gamma, c) \in \Theta$ , a compact subset of  $\mathbb{R}^{2(p+1)+2}$ .

(d) Exponential Smooth Transition AR( $p$ )  $b^\theta(x_1, \dots, x_p) = 1$  and

$$a^\theta(x_1, \dots, x_p) = \mu_1 + \sum_{j=1}^p \phi_{1,j} x_j + \left( \mu_2 + \sum_{j=1}^p \phi_{2,j} x_j \right) (1 - \exp(-\gamma(x_d - c)^2)) \quad (3.15)$$

where  $d \in \{1, \dots, j\}$  and  $\theta = (\mu_i, \phi_{i,j}, i \in \{1, 2\}, j = \{1, \dots, p\}, \gamma, c) \in \Theta$ , a compact subset of  $\mathbb{R}^{2(p+1)+2}$ .

The functional autoregressive models also cover situations for which the regression function is not regular, like the 2 regimes scalar threshold autoregressive model,  $b^\theta(x_1, \dots, x_p) = \sigma_1^2 \mathbb{1}_{\{x_d \leq c\}} + \sigma_2^2 \mathbb{1}_{\{x_d > c\}}$ ,  $d \in \{1, \dots, p\}$  and

$$a^\theta(x_1, \dots, x_p) = \mathbb{1}_{\{x_d \leq c\}} \sum_{j=1}^p \phi_{1,j} x_j + \mathbb{1}_{\{x_d > c\}} \sum_{j=1}^p \phi_{2,j} x_j, \quad (3.16)$$

where  $\theta = (\phi_{i,j}, \sigma_i^2, i \in \{1, 2\}, j \in \{1, \dots, p\}) \in \Theta \subset \mathbb{R}^{2(p+1)}$ .

Assume that for any  $\theta \in \Theta$  and  $(x_1, \dots, x_p) \in \mathbb{R}^p$ ,  $b^\theta(x_1, \dots, x_p) > 0$ . In such a case, the conditional likelihood of the observations may be written as

$$\ln p_\theta(X_{p:n} | X_{0:p-1}) = -\frac{1}{2} \sum_{t=p}^n \ln \left[ 2\pi (b^\theta(X_{t-1}, \dots, X_{t-p}))^2 \right] - \frac{1}{2} \sum_{t=p}^n \frac{\{X_t - a^\theta(X_{t-1}, X_{t-2}, \dots, X_{t-p})\}^2}{(b^\theta(X_{t-1}, X_{t-2}, \dots, X_{t-p}))^2}.$$

The solution of this maximization problem cannot be in general obtained in closed form; a numerical optimization technique is therefore the only option. The choice of “good” initial points is in general mandatory for the optimization algorithm to converge to a sensible optimum. There are several ways to obtain preliminary estimators for these models, but these are in general model dependent.

## 3.2 Consistency and asymptotic normality of the MLE

### 3.2.1 Consistency

We consider first the inference of the parameter  $\theta$  for *misspecified models*. We postulate a model  $\{p^\theta(\cdot) : \theta \in \Theta\}$  for the observations  $(X_0, \dots, X_n)$ . However, the model is misspecified in that the true underlying distribution does not belong to the model. We use the postulated model anyway, and obtain an estimate  $\hat{\theta}_n$  from maximizing the log-likelihood (3.4) where

$$p^\theta(x_{p:n} | x_{0:p-1}) = \prod_{t=p}^n q^\theta(x_{t-p:t-1}; x_t),$$

and  $\{q^\theta : \theta \in \Theta\}$  is a set of Markov kernel densities (with respect to some dominating measure  $\mu$ ) associated to a parametric family of Markov chains of order  $p$ . We derive in this section the asymptotic behavior of  $\hat{\theta}_n$ . Perhaps surprisingly,  $\hat{\theta}_n$  does not behave erratically despite the use of a wrong family of models. First, we show that  $\hat{\theta}_n$  is asymptotically consistent, i.e., converges to a value  $\theta_*$  that maximizes  $\theta \mapsto \mathbb{E}[\ln q^\theta(X_{0:p-1}; X_p)]$  (or a set of value if the solution of this problem is not unique), where the expectation is taken under the true underlying distribution. The density  $q^{\theta_*}$  can be viewed as the “projection” of the true underlying distribution on the model using the Kullback-Leibler divergence, which is defined as  $\mathbb{E}[\ln q^\theta(X_{0:p-1}; X_p)]$ , i.e., as a “distance” measure:  $q^{\theta_*}$  optimizes this quantity over all transition densities in the model. Consider the following assumptions:

**H 3.5** For any  $n \in \mathbb{N}$ , the vector of observations  $(X_0, X_1, \dots, X_n)$  is a realization of a (strict-sense) stationary and ergodic process  $\{X_t, t \in \mathbb{Z}\}$ .

We denote by  $\mathbb{P}$  the probability induced on  $(X^\mathbb{Z}, \mathcal{X}^{\otimes \mathbb{Z}})$  by  $\{X_t, t \in \mathbb{Z}\}$  and by  $\mathbb{E}$  the associated expectation. In particular,  $\{X_t, t \in \mathbb{Z}\}$  is a stationary process but it is not necessarily a stationary Markov chain.

**H 3.6**

- (a)  $\mathbb{P} - \text{a.s.}$ , the function  $\theta \mapsto q^\theta(X_{0:p-1}; X_p)$  is continuous.  
 (b)  $\mathbb{E} [\sup_{\theta \in \Theta} \ln^+ q^\theta(X_{0:p-1}; X_p)] < \infty$

**Theorem 3.7.** Assume **H 3.5** and **H 3.6**. Then, any estimator  $\hat{\theta}_n$  belonging to the set

$$\arg \max_{\theta \in \Theta} \ln p^\theta(X_{p:n} | X_{0:p-1})$$

is strongly consistent in the sense that:

$$\lim_{n \rightarrow \infty} d(\hat{\theta}_n, \Theta_\star) = 0, \mathbb{P} - \text{a.s.} \quad (3.17)$$

where

$$\Theta_\star := \arg \max_{\theta \in \Theta} \mathbb{E} [\ln q^\theta(X_{0:p-1}; X_p)]. \quad (3.18)$$

PROOF. We have for  $n \geq p$ ,

$$\hat{\theta}_n \in \arg \max_{\theta \in \Theta} \bar{L}_n^\theta(X_{p:n} | X_{0:p-1}),$$

where

$$\bar{L}_n^\theta(X_{p:n} | X_{0:p-1}) := (n - p + 1)^{-1} \left( \sum_{t=p}^n \ln q^\theta(X_{t-p:t-1}; X_t) \right).$$

The proof follows from Theorem 3.42 since

- (a)  $\mathbb{E} [\sup_{\theta \in \Theta} \ln^+ q^\theta(X_{0:p-1}; X_p)] < \infty$ ,  
 (b)  $\mathbb{P} - \text{a.s.}$ , the function  $\theta \mapsto \ln q^\theta(X_{0:p-1}; X_p)$  is continuous.

■

Thus, in misspecified models, the MLE strongly converges to the set of parameters that maximize the *relative entropy* (or Kullback-Leibler divergence) between the true distribution and the family of postulated likelihoods. We now consider well-specified models, that is, we assume that  $\{X_t, t \in \mathbb{N}\}$  is the observation process of a Markov chain of order  $p$  associated to the Markov kernel  $Q^\theta$  with  $\theta = \theta_\star \in \Theta$ . In well-specified models, we stress the dependence in  $\theta_\star$  by using the notations

$$\mathbb{P}^{\theta_\star} := \mathbb{P}, \quad \mathbb{E}^{\theta_\star} := \mathbb{E}. \quad (3.19)$$

In this situation, the consistency of the sequence of conditional MLE  $\{\hat{\theta}_n\}_{n \geq 0}$  follows from Theorem 3.7 provided the set  $\Theta_\star$  is reduced to the singleton  $\{\theta_\star\}$ , that is:  $\theta_\star$  is the only parameter  $\theta$  satisfying

$$\mathbb{E}^{\theta_\star} \left[ \ln \frac{q^\theta(X_{0:p-1}; X_p)}{q^{\theta_\star}(X_{0:p-1}; X_p)} \right] = 0.$$

**H 3.8**  $Q^\theta(X_{0:p-1}; \cdot) = Q^{\theta_\star}(X_{0:p-1}; \cdot)$ ,  $\mathbb{P}^{\theta_\star} - \text{a.s.}$  if and only if  $\theta = \theta_\star$ .

The following Corollary is immediate.

**Corollary 3.9** *Under **H** 3.6 and **H** 3.8, if  $\{X_t, t \in \mathbb{N}\}$  is a strict-sense stationary and ergodic sequence under  $\mathbb{P}^{\theta_*}$ , then*

$$\lim_{n \rightarrow \infty} \hat{\theta}_n = \theta_*, \quad \mathbb{P}^{\theta_*} - \text{a.s.}$$

PROOF. According to Theorem 3.7, it is sufficient to show that  $\Theta_* = \{\theta_*\}$ . Now, by the tower property, we have for all  $\theta \in \Theta$ ,

$$\mathbb{E}^{\theta_*} \left[ \ln \frac{q^{\theta_*}(X_{0:p-1}; X_p)}{q^\theta(X_{0:p-1}; X_p)} \right] = \mathbb{E}^{\theta_*} \left[ \mathbb{E}^{\theta_*} \left[ \ln \frac{q^{\theta_*}(X_{0:p-1}; X_p)}{q^\theta(X_{0:p-1}; X_p)} \middle| X_{0:p-1} \right] \right]. \quad (3.20)$$

The RHS is always nonnegative as the expectation of a conditional Kullback-Leibler divergence, which implies that  $\theta_* \in \Theta_*$ . Moreover, using that  $\ln(u) = u - 1$  if and only if  $u = 1$ , (3.20) also shows that if  $\theta \in \Theta_*$ , then

$$q^{\theta_*}(X_{0:p-1}; X_p) = q^\theta(X_{0:p-1}; X_p), \quad \mathbb{P}^{\theta_*} - \text{a.s.}$$

which concludes the proof under **H** 3.8. ■

### 3.2.2 Asymptotic normality

Suppose a sequence of estimators  $\{\hat{\theta}_n\}_{n \geq 0}$  is consistent for a parameter  $\theta_*$  which belongs to an open subset of a Euclidean space. A question of interest relates to the order at which the error  $\hat{\theta}_n - \theta_*$  converges to zero. The answer of course depends on the specificity of the model, but for i.i.d. samples and *regular* statistical models, the order for sensible estimators based on  $n$  observations is  $n^{1/2}$ . Multiplying the estimation error  $\hat{\theta}_n - \theta_*$  by  $n^{1/2}$  produces a proper tradeoff so that, most often, the sequence  $\sqrt{n}(\hat{\theta}_n - \theta_*)$  converges in distribution to a Gaussian random variable with zero-mean and a covariance that can be computed explicitly. This property extends to the Markov chain case without much technical trouble (but of course the expression of the covariance is more involved, except when the model is well-specified). The convergence of  $\sqrt{n}(\hat{\theta}_n - \theta_*)$  is interesting not only from a theoretical point of view but also in practice, since it makes it possible to construct asymptotic confidence regions.

Before going further, just recall briefly the outline of the standard proof in the i.i.d. case. The proof in the Markov case follows almost exactly along the same lines. Let  $X_1, \dots, X_n$  be a sample from some distribution  $\mathbb{P}$ , and let the random and the “true” criterion function be of the form:

$$M_n(\theta) \equiv \frac{1}{n} \sum_{i=1}^n m^\theta(X_i), \quad M(\theta) = \mathbb{E}[m^\theta(X_0)].$$

Assume that the estimator  $\hat{\theta}_n$  is a maximizer of  $M_n(\theta)$  and converges in probability to an element  $\theta_*$  of  $\arg \max_{\theta} M(\theta)$  such that  $\theta_*$  is in the interior of  $\Theta$ . Assume for simplicity that  $\theta$  is one-dimensional. Because  $\hat{\theta}_n \rightarrow_{\mathbb{P}} \theta_*$ , it makes sense to expand  $\dot{M}_n(\hat{\theta}_n)$  in a Taylor series around  $\theta_*$ . Then,

$$0 = \dot{M}_n(\hat{\theta}_n) = \dot{M}_n(\theta_*) + (\hat{\theta}_n - \theta_*) \ddot{M}_n(\tilde{\theta}_n),$$

where  $\tilde{\theta}_n$  is a point between  $\hat{\theta}_n$  and  $\theta_*$ . This can be rewritten as

$$\sqrt{n}(\hat{\theta}_n - \theta_*) = \frac{-\sqrt{n} \dot{M}_n(\theta_*)}{\ddot{M}_n(\tilde{\theta}_n)}. \quad (3.21)$$

If, in this case,  $\mathbb{E}[(\dot{m}^{\theta_*})^2(X_0)]$  is finite, then by the central limit theorem, the numerator  $-\sqrt{n} \dot{M}_n(\theta_*) = -n^{-1/2} \sum_{i=1}^n \dot{m}^{\theta_*}(X_i)$  is asymptotically normal with zero-mean  $\mathbb{E}[\dot{m}^{\theta_*}(X_0)] = \dot{M}(\theta_*) = 0$  and variance  $\mathbb{E}[(\dot{m}^{\theta_*})^2(X_0)]$ .

Next consider the denominator of (3.21). The denominator  $\ddot{M}_n(\tilde{\theta}_n)$  is an average and its limit can be found by using some “uniform” version of the law of large numbers:  $\ddot{M}_n(\tilde{\theta}_n) \rightarrow_{\mathbb{P}} \mathbb{E}[\ddot{m}^{\theta_*}(X_0)]$ , provided

the expectation exists. The difficulty here stems from the fact that  $\tilde{\theta}_n$  is itself a random variable: this is why it is required to use a law of large number that is valid uniformly over at least a vanishing neighborhood of  $\theta_*$ . Together with Slutsky's lemma, these two conclusions yield

$$\sqrt{n}(\hat{\theta}_n - \theta_*) \xrightarrow{\mathcal{L}_P} N\left(0, \frac{\mathbb{E}[(\dot{m}^{\theta_*})^2(X_0)]}{(\mathbb{E}[\dot{m}^{\theta_*}(X_0)])^2}\right). \quad (3.22)$$

The preceding derivation can be made rigorous by imposing appropriate technical conditions (on the regularity of the function  $\theta \mapsto \dot{M}_n(\theta)$ , the validity of interchanging derivation and expectation), often called *regularity conditions* in the i.i.d. case. Perhaps surprisingly, the most challenging part of the proof consists of showing that  $\dot{M}_n(\tilde{\theta}_n)$  is converging to  $\mathbb{E}[\dot{m}^{\theta_*}(X_0)]$  (see above); all the other steps are direct consequences of classical limit theorems.

The derivation can be extended to higher-dimensional parameters. For a  $d$ -dimensional parameter, we use  $d$  estimating equations  $\dot{M}_n(\theta) = 0$ , where  $\dot{M}_n : \mathbb{R}^d \mapsto \mathbb{R}^d$ . The derivatives  $\dot{M}_n(\theta_*)$  are  $(d \times d)$ -matrices that converge to the  $(d \times d)$  matrix  $\Gamma_* := \mathbb{E}[\dot{m}^{\theta_*}(X_0)]$  with entries  $\mathbb{E}[\partial^2 m^{\theta_*}(X_0)/\partial \theta_i \partial \theta_j]$ . The limiting distribution (see (3.22)) may be expressed as

$$\sqrt{n}(\hat{\theta}_n - \theta_*) \xrightarrow{\mathcal{L}_P} N_d(0, \Gamma_*^{-1} \Sigma_* \Gamma_*^{-1}), \quad (3.23)$$

where  $\Sigma_* = \mathbb{E}[\dot{m}^{\theta_*}(X_0)(\dot{m}^{\theta_*}(X_0))']$  is the covariance matrix of the “pseudo-scores.” Note that the matrix  $\Gamma_*$  should be nonsingular.

In this section we extend these results to Markov chains. We consider first the case of (possibly) misspecified models. Assume that  $\Theta \subset \mathbb{R}^d$  and recall that  $\Theta^o$  is the interior of  $\Theta$ .

**H 3.10** *There exists  $\theta_* \in \Theta^o$  such that  $\hat{\theta}_n \xrightarrow{\mathbb{P}\text{-prob}} \theta_*$ .*

### H 3.11

- (a) *The function  $\theta \mapsto q^\theta(X_{0:p-1}; X_p)$  is  $\mathbb{P}$ -a.s. twice continuously differentiable in an open neighborhood of  $\theta_*$ .*
- (b) *There exists  $\rho \in (0, 1)$  such that for all indexes  $i, j \in \{1, \dots, d\}$ ,*

$$\mathbb{E} \left[ \sup_{\theta \in B(\theta_*, \rho)} \left| \frac{\partial^2 \ln q^\theta}{\partial \theta_i \partial \theta_j}(X_{0:p-1}; X_p) \right| \right] < \infty.$$

**H 3.12** *There exists a  $d \times d$  nonsingular matrix  $\Sigma_*$  such that*

$$n^{-1/2} \sum_{t=p}^n (\nabla \ln q^{\theta_*}(X_{t-p:t-1}; X_t) - \mathbb{E}[\nabla \ln q^{\theta_*}(X_{0:p-1}; X_p)]) \xrightarrow{\mathcal{L}_P} N_d(0, \Sigma_*). \quad (3.24)$$

**Theorem 3.13.** *Assume H 3.10, H 3.11 and H 3.12. In addition, assume that  $\{X_t, t \in \mathbb{N}\}$  is a strict-sense stationary and ergodic sequence under  $\mathbb{P}$  and that  $\Gamma_* := \mathbb{E}[\nabla^2 \ln q^{\theta_*}(X_{0:p-1}; X_p)]$  is nonsingular. Then,*

$$n^{1/2}(\hat{\theta}_n - \theta_*) \xrightarrow{\mathcal{L}_P} N(0, \Gamma_*^{-1} \Sigma_* \Gamma_*^{-1})$$



where  $\Sigma_\star$  is defined in (3.24).

We preface the proof by the following technical Lemma, which will be useful in misspecified and well-specified models as well.

**Lemma 3.14** Assume that **H 3.11** hold for some  $\theta_\star \in \Theta$ . Assume in addition that  $\{X_t, t \in \mathbb{N}\}$  is a strict-sense stationary and ergodic sequence under  $\mathbb{P}$  and let  $\{\theta_n, n \in \mathbb{N}\}$  be a sequence of random vectors such that  $\theta_n \rightarrow_{\mathbb{P}} \theta_\star$ . Then, for all  $i, j \in \{1, \dots, d\}$ ,

$$n^{-1} \sum_{t=p}^n \frac{\partial^2 \ln q^{\theta_n}(X_{t-p:t-1}; X_t)}{\partial \theta_i \partial \theta_j} \xrightarrow{\mathbb{P}\text{-prob}} \mathbb{E} \left[ \frac{\partial^2 \ln q^{\theta_\star}(X_{0:p-1}; X_p)}{\partial \theta_i \partial \theta_j} \right].$$

PROOF. Denote  $A_t(\theta) = \frac{\partial^2 \ln q^\theta(X_{t-p:t-1}; X_t)}{\partial \theta_i \partial \theta_j}$ . Since by the Birkhoff ergodic theorem,  $n^{-1} \sum_{t=p}^n A_t(\theta_\star) \rightarrow_{\mathbb{P}} \mathbb{E}[A(\theta_\star)]$ , we only need to show that

$$n^{-1} \sum_{t=p}^n |A_t(\theta_\star) - A_t(\theta_n)| \xrightarrow{\mathbb{P}\text{-prob}} 0. \quad (3.25)$$

Let  $\varepsilon > 0$  and choose  $0 < \eta < \rho$  such that

$$\mathbb{E} \left( \sup_{\theta \in B(\theta_\star, \eta)} |A_p(\theta_\star) - A_p(\theta)| \right) < \varepsilon. \quad (3.26)$$

The existence of such  $\eta$  follows from the  $\mathbb{P}$ -a.s. continuity of  $\theta \mapsto A_t(\theta)$  under **H 3.11**-(a)) and by the Lebesgue convergence theorem under **H 3.11**-(b)). We then have

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{P} \left( n^{-1} \sum_{t=p}^n |A_t(\theta_\star) - A_t(\theta_n)| \geq \varepsilon, \theta_n \in B(\theta_\star, \eta) \right) \\ \leq \limsup_{n \rightarrow \infty} \mathbb{P} \left( n^{-1} \sum_{t=p}^n \sup_{\theta \in B(\theta_\star, \eta)} |A_t(\theta_\star) - A_t(\theta)| \geq \varepsilon \right) = 0, \end{aligned}$$

where the last equality follows from (3.26) and the Birkhoff ergodic theorem. Moreover, since  $\hat{\theta}_n \xrightarrow{\mathbb{P}\text{-prob}} \theta_\star$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(\theta_n \notin B(\theta_\star, \eta)) = 0$  so that finally,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left( n^{-1} \sum_{t=p}^n |A_t(\theta_\star) - A_t(\theta_n)| \geq \varepsilon \right) = 0.$$

Thus, (3.25) holds and the proof follows.  $\blacksquare$

PROOF. [of Theorem 3.13] Since  $\hat{\theta}_n$  cancels the derivatives of the log-likelihood, a Taylor expansion at the point  $\theta = \theta_\star$  with an integral form of the remainder yields

$$\begin{aligned} n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\hat{\theta}_n}(X_{t-p:t-1}; X_t) &= 0 = n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\theta_\star}(X_{t-p:t-1}; X_t) \\ &\quad + n^{-1} \sum_{t=p}^n \left( \int_0^1 \nabla^2 \ln q^{\theta_{n,s}}(X_{t-p:t-1}; X_t) ds \right) \sqrt{n}(\hat{\theta}_n - \theta_\star), \end{aligned} \quad (3.27)$$

where  $\theta_{n,s} = s\hat{\theta}_n + (1-s)\theta_\star$ . According to **H 3.12**, we have

$$n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\theta_\star}(X_{t-p:t-1}; X_t) \xrightarrow{\mathcal{L}_{\theta_\star}} \mathcal{N}(0, \Sigma_\star).$$

We now turn to the asymptotic properties of the second term of the right-hand side in (3.27). More precisely, since by **H 3.10**,  $\hat{\theta}_n \rightarrow_{\mathbb{P}} \theta_\star$ , Lemma 3.14 shows that for all  $i, j \in \{1, \dots, d\}$ ,

$$n^{-1} \sum_{t=p}^n \int_0^1 \frac{\partial^2 \ln q^{\theta_{n,s}}}{\partial \theta_i \partial \theta_j}(X_{t-p:t-1}; X_t) ds \xrightarrow{\mathbb{P}\text{-prob}} \mathbb{E}^{\theta_\star} \left[ \frac{\partial^2 \ln q^{\theta_\star}}{\partial \theta_i \partial \theta_j}(X_{0:p-1}; X_p) \right].$$

This implies that

$$n^{-1} \sum_{t=p}^n \left( \int_0^1 \nabla^2 \ln q^{\theta_{n,s}}(X_{t-p:t-1}; X_t) ds \right) \xrightarrow{\mathbb{P}^{\theta_*} - \text{prob}} \Gamma_* := \mathbb{E}[\nabla^2 \ln q^{\theta_*}(X_{0:p-1}; X_p)] .$$

The proof then follows by applying the Slutsky Lemma to (3.27).  $\blacksquare$

We now turn to well-specified models, that is, the observation process  $\{X_n, n \in \mathbb{N}\}$  is the realization of a Markov chain associated to an unknown parameter  $\theta_*$ .

Before stating the result, we need some additional assumptions. The parameter set  $\Theta$  is now a compact subset of  $\mathbb{R}^d$  with a non-empty interior and we assume that  $\theta_* \in \Theta^o$ , where  $\Theta^o$  is the interior of  $\Theta$ . Moreover,

### H 3.15

- (a) The function  $\theta \mapsto q^\theta(X_{0:p-1}; X_p)$  is,  $\mathbb{P}^{\theta_*}$  - a.s., twice continuously differentiable in an open neighborhood of  $\theta_*$ ,
- (b) there exists  $\rho > 0$  such that for all  $i, j \in \{1, \dots, d\}$ ,

$$\mathbb{E}^{\theta_*} \left[ \sup_{\theta \in B(\theta_*, \rho)} \left| \frac{\partial^2 \ln q^\theta}{\partial \theta_i \partial \theta_j}(X_{0:p-1}; X_p) \right| \right] < \infty .$$

H 3.15 is similar to H 3.11. In well-specified models, it is possible to obtain the weak convergence of the score function using a martingale argument (see Exercise 3.43). Some additional assumptions are needed.

**H 3.16** *There exists  $\rho > 0$  such that*

$$\mathbb{E}^{\theta_*} \left[ \left| \frac{\partial \ln q^{\theta_*}}{\partial \theta_i}(X_{0:p-1}; X_p) \right|^2 \right] < \infty , \quad (3.28)$$

$$\int \sup_{\theta \in B(\theta_*, \rho)} \left| \frac{\partial q^\theta}{\partial \theta_i}(X_{0:p-1}; x_p) \right| d\mu(x_p) < \infty , \quad \mathbb{P}^{\theta_*} - \text{a.s.} , \quad (3.29)$$

$$\int \sup_{\theta \in B(\theta_*, \rho)} \left| \frac{\partial^2 q^\theta}{\partial \theta_i \partial \theta_j}(X_{0:p-1}; x_p) \right| d\mu(x_p) < \infty , \quad \mathbb{P}^{\theta_*} - \text{a.s.} , \quad (3.30)$$

and the Fisher information matrix  $\mathcal{J}(\theta_*)$  defined by

$$\begin{aligned} \mathcal{J}(\theta_*) &:= -\mathbb{E}^{\theta_*} \left[ \nabla^2 \ln q^{\theta_*}(X_{0:p-1}; X_p) \right] \\ &= \mathbb{E}^{\theta_*} \left[ \nabla \ln q^{\theta_*}(X_{0:p-1}; X_p) \nabla \ln q^{\theta_*}(X_{0:p-1}; X_p)^T \right] \end{aligned} \quad (3.31)$$

is nonsingular.

Note that the Fisher information matrix has two expressions given in (3.31). These two expressions are classically obtained under the assumption (3.30). We now state and prove the weak convergence of the score function using a martingale-type approach.

**Lemma 3.17** *Under H 3.16,*

$$n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\theta_*}(X_{t-p:t-1}; X_t) \xrightarrow{\mathcal{L}_{\mathbb{P}^{\theta_*}}} N(0, \mathcal{J}(\theta_*)) ,$$

where  $\mathcal{J}(\theta_*)$  is defined in (3.31).

PROOF. Let  $\mathcal{F}$  be the filtration  $\mathcal{F} = (\mathcal{F}_n)_{n \in \mathbb{N}}$  where  $\mathcal{F}_n = \sigma(X_0, \dots, X_n)$  and let

$$M_n := \sum_{t=p}^n \nabla \ln q^{\theta_*}(X_{t-p:t-1}; X_t) .$$

Note that according to (3.28),  $\mathbb{E}(|M_t|^2) < \infty$  and that

$$\begin{aligned} \mathbb{E}^{\theta_*} \left[ \nabla \ln q^{\theta_*}(X_{t-p:t-1}; X_t) \middle| \mathcal{F}_{t-1} \right] &= \int \nabla q^{\theta_*}(X_{t-p:t-1}; x_t) dx_t \\ &= \nabla \int q^{\theta_*}(X_{t-p:t-1}; x_t) dx_t \Big|_{\theta=\theta_*} = 0 , \end{aligned}$$

where the assumption (3.29) allows us to interchange  $\int$  and  $\nabla$ . Finally,  $\{M_t, t \geq p\}$  is a square integrable  $\mathcal{F}$ -martingale with stationary increments. The proof follows by applying Theorem A.7. ■

We now have all the ingredients for obtaining the central limit theorem of the MLE in well-specified models.

**Theorem 3.18.** Assume that  $\hat{\theta}_n \xrightarrow{\mathbb{P}^{\theta_*-prob}} \theta_*$  and that **H 3.15** and **H 3.16** hold. Then,

$$\sqrt{n}(\hat{\theta}_n - \theta_*) \xrightarrow{\mathcal{L}_{\mathbb{P}}} N(0, \mathcal{J}^{-1}(\theta_*)) ,$$

where  $\mathcal{J}(\theta_*)$  is defined in (3.31).

PROOF. See Exercise 3.48 ■

Note that  $\mathcal{J}(\theta_*)$  is the asymptotic Fisher information matrix in the sense that

$$\mathcal{J}(\theta_*) = - \lim_{n \rightarrow \infty} n^{-1} \mathbb{E}^{\theta_*} \left[ \nabla^2 \ln \left( \prod_{k=p}^n q^{\theta_*}(X_{k-p:k-1}; X_k) \right) \middle|_{\theta=\theta_*} \right] .$$

**Example 3.19 (Autoregressive models)** Let  $(\phi_{*,1}, \phi_{*,2}, \dots, \phi_{*,p}) \in \mathbb{R}^p$  be coefficients such that  $\phi_*(z) = 1 - \sum_{j=1}^p \phi_{*,j} z^j \neq 0$  for  $|z| \leq 1$ , where  $p > 0$  is an integer. Assume that  $\{X_t, t \in \mathbb{Z}\}$  is the causal autoregressive process:  $X_t = \sum_{j=1}^p \phi_{*,j} X_{t-j} + \sigma_* Z_t$  where  $\{Z_t, t \in \mathbb{Z}\}$  is white Gaussian noise with unit-variance. Set  $\theta = (\phi_1, \dots, \phi_p, \sigma^2)$ ,  $\theta_* = (\phi_{*,1}, \dots, \phi_{*,p}, \sigma_*^2)$  and

$$\ln q^{\theta}(x_{0:p-1}, x_p) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \left( x_p - \sum_{j=1}^p \phi_j x_{p-j} \right)^2 .$$

By differentiating twice with respect to  $\sigma^2$  and  $\phi_i$ , we get

$$\begin{aligned} \frac{\partial^2 \ln q^{\theta}(X_{0:p-1}; X_p)}{\partial \phi_i \partial \phi_j} &= -\frac{1}{\sigma^2} X_{p-i} X_{p-j} , & (i, j) \in \{1, \dots, p\}^2 \\ \frac{\partial^2 \ln q^{\theta}(X_{0:p-1}; X_p)}{\partial \phi_j \partial \sigma^2} &= -\frac{1}{\sigma^4} \left( X_p - \sum_{j=1}^p \phi_j X_{p-j} \right) X_{p-j} , & j \in \{1, \dots, p\} \\ \frac{\partial^2 \ln q^{\theta}(X_{0:p-1}; X_p)}{\partial \sigma^2 \partial \sigma^2} &= \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} \left( X_p - \sum_{j=1}^p \phi_j X_{p-j} \right)^2 . \end{aligned}$$

Since, for any  $i \in \{1, \dots, p\}$ ,

$$\mathbb{E}^{\theta_*} \left[ \left( X_p - \sum_{j=1}^p \phi_{*,j} X_{p-j} \right) X_{p-i} \right] = \mathbb{E}^{\theta_*} [Z_p X_{p-i}] = 0,$$

and, similarly

$$\mathbb{E}^{\theta_*} \left[ \left( X_p - \sum_{j=1}^p \phi_{*,j} X_{p-j} \right)^2 \right] = \mathbb{E}^{\theta_*} [Z_p^2] = \sigma_*^2,$$

we get that  $\mathcal{J}(\theta_*)$  is a block diagonal matrix, with block diagonal elements equal to  $1/2\sigma_*^4$  and  $\sigma_*^{-2}\Gamma_*$ , where  $\Gamma_*$  is the  $p \times p$  covariance matrix,  $[\Gamma_*]_{i,j} = \text{Cov}^{\theta_*}(X_0, X_{i-j})$ ,  $1 \leq i, j \leq p$ . If  $p = 1$ , then  $\Gamma_* = \sigma_*^2/(1 - \phi_{*,1}^2)$ , and the asymptotic variance for the autoregressive parameter is equal to  $1/(1 - \phi_{*,1}^2)$ . For  $p = 2$ , the reader can verify (Exercise 3.46) that

$$\sigma_*^{-2}\Gamma_* = \begin{bmatrix} 1 - \phi_{*,2}^2 & -\phi_{*,1}(1 + \phi_{*,2}) \\ -\phi_{*,1}(1 + \phi_{*,2}) & 1 - \phi_{*,2}^2 \end{bmatrix}. \quad (3.32)$$

Note that the asymptotic variance of the estimate of  $\phi_{*,1}$  depends only on  $\phi_{*,2}$ .

**Example 3.20 (Threshold autoregressive models)** We use the notations and definitions of Example 3.3. Denote by  $\theta = [\eta_1, \eta_2, r, d]$  with  $\eta_i = [(\phi_{i,j})_{j=0}^{p_i}, \sigma_i^2]$ ,  $i = 1, 2$ , the unknown parameters. In this case the log-likelihood may be expressed as

$$\begin{aligned} \ln q^\theta(x_{0:p-1}; x_p) &= \left\{ -\frac{1}{2} \ln(2\pi\sigma_1^2) - \frac{1}{2\sigma_1^2} \left( X_p - \phi_{1,0} - \sum_{j=1}^{p_1} \phi_{1,j} X_{p-j} \right)^2 \right\} \mathbb{1}_{\{X_{t-d} \leq r\}} \\ &\quad + \left\{ -\frac{1}{2} \ln(2\pi\sigma_2^2) - \frac{1}{2\sigma_2^2} \left( X_p - \phi_{2,0} - \sum_{j=1}^{p_2} \phi_{2,j} X_{p-j} \right)^2 \right\} \mathbb{1}_{\{X_{t-d} > r\}}, \end{aligned}$$

where  $p = \max(p_1, p_2)$ . When the delay  $d$  and the threshold  $r$  are known, Theorem 3.7 shows that, provided that  $\{X_t, t \in \mathbb{Z}\}$  is an ergodic sequence, then the regression parameters  $\hat{\eta}_i(r, d) = [\hat{\phi}_{i,j}(r, d), \hat{\sigma}_i^2(r, d)]_{j=1}^{p_i}$ ,  $i = 1, 2$  are consistent in the sense of (3.17). It has been established by ?, Theorem 1 that the same remains true if the delay and the threshold are estimated as the maximum of the profile likelihood.

Assume now that  $\{X_t, t \in \mathbb{Z}\}$  is a strict sense stationary geometrically ergodic SETAR process

$$\begin{aligned} X_t &= \left\{ \phi_{*,1,0} + \sum_{j=1}^{p_1} \phi_{*,1,j} X_{t-j} + \sigma_{*,1} Z_t \right\} \mathbb{1}_{\{X_{t-d_*} \leq r_*\}} \\ &\quad + \left\{ \phi_{*,2,0} + \sum_{j=1}^{p_2} \phi_{*,2,j} X_{t-j} + \sigma_{*,2} Z_t \right\} \mathbb{1}_{\{X_{t-d_*} > r_*\}}, \end{aligned} \quad (3.33)$$

where  $\{Z_t, t \in \mathbb{Z}\}$  is a strong white Gaussian noise with zero-mean and unit-variance. Such condition is satisfied, for example, if  $\max_{i=1,2} \sum_{j=1}^{p_i} |\phi_{*,i,j}| < 1$  (see ??, condition (??)). Assume first that the delay  $d$  and the threshold  $r$  are both known:  $r = r_*$  and  $d = d_*$ . We may then apply Theorem 3.18 to show that the distribution of the parameters in the two regimes are asymptotically normal,

$$\sqrt{n}(\hat{\eta}_1(r_*, d_*) - \eta_{*,1}, \hat{\eta}_2(r_*, d_*) - \eta_{*,2}) \xrightarrow{\mathcal{L}_P} \mathbf{N}(0, \mathcal{J}^{-1}(\theta_*))$$

where  $\theta_* = [\eta_{*,1}, \eta_{*,2}, r_*, d_*]$ ,  $\eta_{*,i} = [(\phi_{*,i,j})_{j=0}^{p_i}, \sigma_{*,i}^2]$ ,  $i = 1, 2$ . Note that, for any  $(i, j) \in \{1, 2\}$ ,  $i \neq j$ , and any  $(k, \ell) \in \{1, \dots, p_i\}^2$ , we get

$$\frac{\partial^2 \ln q^\theta(x_{0:p-1}; x_p)}{\partial \phi_{i,k} \partial \phi_{j,\ell}} = 0, \quad \frac{\partial^2 \ln q^\theta(x_{0:p-1}; x_p)}{\partial \sigma_i^2 \partial \phi_{j,k}} = 0,$$

showing that the matrix  $\mathcal{J}(\theta_*)$  is block diagonal and therefore that the estimators in the two regimes are asymptotically independent. The entries of the Fisher information matrix are given by

$$\begin{aligned}\mathbb{E}^{\theta_*} \left[ \frac{\partial^2 \ln q^{\theta_*}(X_{0:p-1}; X_p)}{\partial \phi_{1,k} \partial \phi_{1,\ell}} \right] &= \frac{1}{\sigma_1^2} \mathbb{E}^{\theta_*} [X_{p-k} X_{p-\ell} \mathbb{1}_{\{X_{p-d} \leq r\}}], & (k, \ell) \in \{1, \dots, p_1\} \\ \mathbb{E}^{\theta_*} \left[ \frac{\partial^2 \ln q^{\theta_*}(X_{0:p-1}; X_p)}{\partial \phi_{1,0} \partial \phi_{1,k}} \right] &= \frac{1}{\sigma_1^2} \mathbb{E}^{\theta_*} [X_{p-k} \mathbb{1}_{\{X_{p-d} \leq r\}}], & k \in \{1, \dots, p_1\} \\ \mathbb{E}^{\theta_*} \left[ \frac{\partial^2 \ln q^{\theta_*}(X_{0:p-1}; X_p)}{\partial \phi_{1,k} \partial \sigma_1^2} \right] &= 0, \\ \mathbb{E}^{\theta_*} \left[ \frac{\partial^2 \ln q^{\theta_*}(X_{0:p-1}; X_p)}{\partial \sigma_1^2 \partial \sigma_1^2} \right] &= \frac{1}{2\sigma_1^2} \mathbb{P}^{\theta_*}(X_0 \leq r).\end{aligned}$$

The Fisher information matrix for the parameters of the upper regime can be written similarly. When the threshold  $r$  and  $d$  are both unknown, the assumptions of Theorem 3.18 are no longer satisfied. It is shown in ?, Theorem 2 that

$$\{n(\hat{r}_n - r_*), \sqrt{n}(\hat{\eta}_1 - \eta_{*,1}, \hat{\eta}_2 - \eta_{*,2})\}$$

converges in distribution, that  $n(\hat{r}_n - r_*)$  and  $\sqrt{n}(\hat{\eta}_1 - \eta_{*,1}, \hat{\eta}_2 - \eta_{*,2})$  are asymptotically independent. ? identified the distribution of  $n(\hat{r}_n - r_*)$ : this is an interval on which a compound Poisson process (whose parameters depend on  $\theta_*$ ) attains its global minimum.

In general, the problem of determining the set of parameters  $\Theta_*$  for  $\{X_t, t \in \mathbb{Z}\}$  to be geometrically ergodic is still an open problem. For  $p_1 = p_2 = 1$  and  $d = 1$ , however, the problem is solved. The condition in this case is  $\phi_{*,1,1} < 1$ ,  $\phi_{*,2,1} < 1$  and  $\phi_{*,1,1}\phi_{*,2,1} < 1$ . For details, see ?.

### 3.3 Observation-driven models

Let  $(X, d)$  and  $(Y, \delta)$  be Polish spaces equipped with their Borel sigma-field  $\mathcal{X}$  and  $\mathcal{Y}$ . Consider  $\{Q^\theta, \theta \in \Theta\}$  a family of Markov kernels on  $X \times \mathcal{Y}$  indexed by  $\theta \in \Theta$  where  $(\Theta, d)$  is a compact metric space. Assume that all  $(\theta, x) \in \Theta \times X$ ,  $Q^\theta(x, \cdot)$  is dominated by some  $\sigma$ -finite measure  $\mu$  on  $(Y, \mathcal{Y})$  and denote by  $q^\theta(x, \cdot)$  its Radon-Nikodym derivative:  $q^\theta(x, y) = dQ^\theta(x, \cdot)/d\mu(y)$ . Finally, let  $\{f^\theta : \theta \in \Theta\}$  be a family of measurable functions from  $(X \times Y, \mathcal{X} \otimes \mathcal{Y})$  to  $(X, \mathcal{X})$ . We consider the following observation-driven model (see Theorem 2.21)

$$\begin{aligned}\mathbb{P}^\theta[Y_t \in A | \mathcal{F}_{t-1}] &= Q^\theta(X_{t-1}, A) = \int_A q^\theta(X_{t-1}, y) \mu(dy), \quad \text{for any } A \in \mathcal{Y}, \\ X_t &= f_{Y_t}^\theta(X_{t-1}).\end{aligned}$$

With these notations, the distribution of  $(Y_1, \dots, Y_n)$  conditionally on  $X_0 = x$  has a density with respect to the product measure  $\mu^{\otimes n}$  given by

$$y_{1:n} \mapsto \prod_{t=1}^n q^\theta(f_{y_{1:t-1}}^\theta(x), y_t), \quad (3.34)$$

where we have set for all  $s \leq t$  and all  $y_{s:t} \in Y^{t-s+1}$ ,

$$f_{y_{s:t}}^\theta = f_{y_t}^\theta \circ f_{y_{t-1}}^\theta \circ \dots \circ f_{y_s}^\theta, \quad (3.35)$$

with the convention  $f_{y_{1:0}}^\theta(x_0) = x_0$ . Note that, for any  $t \geq 0$ ,  $X_t$  is a deterministic function of  $Y_{1:t}$  and  $X_0$ , i.e.,

$$X_t = f_{Y_{1:t}}^\theta(X_0) = f_{Y_t}^\theta \circ f_{Y_{t-1}}^\theta \circ \dots \circ f_{Y_1}^\theta(X_0). \quad (3.36)$$

In this section, we study the asymptotic properties of  $\hat{\theta}_{n,x}$ , the conditional Maximum Likelihood Estimator (MLE) of the parameter  $\theta$  based on the observations  $(Y_1, \dots, Y_n)$  and associated to the parametric family of likelihood functions given in (3.34), that is, we consider

$$\hat{\theta}_{n,x} \in \arg \max_{\theta \in \Theta} L_{n,x}^{\theta}(Y_{1:n}), \quad (3.37)$$

where

$$L_{n,x}^{\theta}(Y_{1:n}) := n^{-1} \ln \left( \prod_{t=1}^n q^{\theta}(f_{Y_{1:t-1}}^{\theta}(x), Y_t) \right). \quad (3.38)$$

We are especially interested here in inference for *misspecified models*, that is, we *do not assume* that the distribution of the observations belongs to the set of distributions where the maximization occurs. In particular,  $\{Y_t, t \in \mathbb{Z}\}$  are not necessarily the observation process associated to the recursion (3.36).

Nevertheless, consider the following assumptions

**H 3.21**  $\{Y\}$  is a strict-sense stationary and ergodic stochastic process.

Under **H 3.21**, denote by  $\mathbb{P}$  the distribution of  $\{Y\}$  on  $(Y^{\mathbb{Z}}, \mathcal{Y}^{\mathbb{Z}})$ . Write  $\mathbb{E}$ , the associated expectation.

**H 3.22** For all  $(x, y) \in X \times Y$ , the functions  $\theta \mapsto f_y^{\theta}(x)$  and  $(y, \theta) \mapsto q^{\theta}(x, y)$  are continuous.

**H 3.23** There exists a family of  $\mathbb{P}$  – a.s. finite random variables

$$\left\{ f_{Y_{-\infty:t}}^{\theta} : (\theta, t) \in \Theta \times \mathbb{Z} \right\}$$

such that for all  $x \in X$ ,

- (a)  $\lim_{m \rightarrow \infty} \sup_{\theta \in \Theta} d(f_{Y_{-m:0}}^{\theta}(x), f_{Y_{-\infty:0}}^{\theta}) = 0$ ,  $\mathbb{P}$  – a.s.,
- (b)  $\lim_{t \rightarrow \infty} \sup_{\theta \in \Theta} |\ln q^{\theta}(f_{Y_{1:t-1}}^{\theta}(x), Y_t) - \ln q^{\theta}(f_{Y_{-\infty:t-1}}^{\theta}, Y_t)| = 0$ ,  $\mathbb{P}$  – a.s.,
- (c)  $\mathbb{E} \left[ \sup_{\theta \in \Theta} \left( \ln q^{\theta}(f_{Y_{-\infty:t-1}}^{\theta}, Y_t) \right)_+ \right] < \infty$ .

In the following, we set for all  $(\theta, t) \in \Theta \times \mathbb{N}$ ,

$$\bar{\ell}^{\theta}(Y_{-\infty:t}) := \ln q^{\theta}(f_{Y_{-\infty:t-1}}^{\theta}, Y_t). \quad (3.39)$$

**Remark 3.24** When checking **H 3.23**, we usually introduce  $f_{Y_{-\infty:0}}^{\theta}$  by showing that for all  $(\theta, x) \in \Theta \times X$ ,  $f_{Y_{-m:0}}^{\theta}(x)$  converges,  $\mathbb{P}$  – a.s., as  $m$  goes to infinity to a limit that does not depend on  $x$ . We can therefore denote by  $f_{Y_{-\infty:0}}^{\theta}$  this limit. With this definition, we then check **H 3.23**–((a))–((b))–((c)).

Note that under **H 3.22**,  $\hat{\theta}_{n,x}$  is well-defined. The following theorem establishes the consistency of the sequence of estimators  $\{\hat{\theta}_{n,x}, n \in \mathbb{N}\}$ .

**Theorem 3.25.** Assume **H 3.21**, **H 3.22** and **H 3.23**. Then, for all  $x \in X$ ,

$$\lim_{n \rightarrow \infty} d(\hat{\theta}_{n,x}, \Theta_*) = 0, \quad \mathbb{P} - \text{a.s.}$$

where  $\Theta_* := \arg \max_{\theta \in \Theta} \mathbb{E}[\bar{\ell}^\theta(Y_{-\infty;1})]$  and  $\bar{\ell}^\theta(Y_{-\infty;1})$  is defined in (3.39).

PROOF. The proof directly follows from Theorem 3.42 provided that

$$\mathbb{E}[\sup_{\theta \in \Theta} (\bar{\ell}^\theta(Y_{-\infty;1}))_+] < \infty, \quad (3.40)$$

$$\text{the function } \theta \mapsto \bar{\ell}^\theta(Y_{-\infty;1}) \text{ is upper-semicontinuous, } \mathbb{P} - \text{a.s.}, \quad (3.41)$$

$$\lim_{n \rightarrow \infty} \sup_{\theta \in \Theta} |\mathbb{L}_{n,x}^\theta(Y_{1:n}) - \bar{\mathbb{L}}_n^\theta(Y_{-\infty;n})| = 0, \quad \mathbb{P} - \text{a.s.}, \quad (3.42)$$

where  $\bar{\mathbb{L}}_n^\theta(Y_{-\infty;n}) = n^{-1} \sum_{t=1}^n \bar{\ell}^\theta(Y_{-\infty;t})$ .

But, (3.40) follows from **H** 3.23-(c), (3.41) follows by combining **H** 3.23-(a) and **H** 3.22 since a uniform limit of continuous functions is continuous and (3.42) is direct from **H** 3.23-(b) and the definitions of  $\mathbb{L}_{n,x}^\theta(Y_{1:n})$  and  $\bar{\mathbb{L}}_n^\theta(Y_{-\infty;n})$ . The proof is completed.  $\blacksquare$

**Example 3.26 (GARCH(1, 1))** Assume that the observations  $\{Y_t, t \in \mathbb{Z}\}$  are a strict-sense stationary ergodic process. We fit to the observations a GARCH(1, 1) model with Student  $t$ -innovations,

$$\begin{cases} Y_t = \sigma_t \left( \frac{v-2}{v} \right)^{1/2} \varepsilon_t \\ \sigma_t^2 = \alpha_0 + \alpha_1 Y_{t-1}^2 + \beta_1 \sigma_{t-1}^2, \quad t \in \mathbb{Z} \end{cases} \quad (3.43)$$

where  $\{\varepsilon_t, t \in \mathbb{N}\}$  is an i.i.d. sequence of  $t_v$ -distributed random variables (where  $v$  is the number of degrees of freedom) and  $\theta = (v, \alpha_0, \alpha_1, \beta_1) \in \Theta$  is a compact subset of

$$\{(v, \alpha_0, \alpha_1, \beta_1), v > 2, \alpha_0 > 0, \alpha_1 > 0, \beta_1 > 0, \alpha_1 + \beta_1 < 1\}.$$

The restriction on the degrees of freedom parameter  $v$  ensures the conditional variance to be finite. The variance of  $\varepsilon_1$  is equal to  $v/(v-2)$ , therefore,  $((v-2)/v)^{1/2} \varepsilon_1$  has unit-variance. Recall that the density of a Student  $t$ -distribution with  $v$  degrees of freedom is given by

$$t_v(z) = \frac{\Gamma(\frac{v+1}{2})}{\Gamma(\frac{v}{2})\Gamma(\frac{v}{2})} v^{-1/2} [1 + z^2/v]^{-(v+1)/2}.$$

Setting  $X_t = \sigma_{t+1}^2$ , (3.43) defines an observation-driven model with  $q^\theta(x, y) = t_v(y/\sqrt{x})$  and  $f_y^\theta(x) = \alpha_0 + \alpha_1 y^2 + \beta_1 x$ . Given initial values  $Y_0$  and  $\sigma_0^2$  to be specified below, the conditional log-likelihood may be expressed as

$$\mathbb{L}_{n, \sigma_0}^\theta(Y_{1:n}) = \sum_{t=1}^n \log t_v(Y_t / \sigma_t),$$

where  $\sigma_t^2$  are computed recursively using (3.43). For a given value of  $\theta$ , the unconditional variance (corresponding to the stationary value of the variance) is a sensible choice for the unknown initial values  $\sigma_0^2$

$$\sigma_0^2 = \frac{\alpha_0}{1 - \alpha_1 - \beta_1}. \quad (3.44)$$

In this case, for any integer  $m$ , we have (see (2.21))

$$f_{Y_{-m:0}}^\theta(x) = \beta_1^{m+1} x + \sum_{j=0}^m \beta_1^j (\alpha_0 + \alpha_1 Y_{-j}^2).$$

Assume that  $\mathbb{E}_*[Y_0^2] < \infty$ . Since  $\Theta$  is compact, there exist  $b_1 < 1$ ,  $a_0 > 0$  and  $a_1 < 1$  such that, for all  $\theta \in \Theta$ ,  $0 \leq \beta_1 \leq b_1$ ,  $\alpha_0 \leq a_0$  and  $\alpha_1 \leq a_1$ . The series  $\sum b_1^j (a_0 + a_1 Y_{-j}^2)$  converges  $\mathbb{P}_*$ -a.s. Define

$$f_{Y_{-\infty};0}^\theta = \sum_{j=0}^{\infty} \beta_1^j (\alpha_0 + \alpha_1 Y_{-j}^2).$$

With these definitions, we get

$$\sup_{\theta \in \Theta} |f_{Y_{-m};0}^\theta(x) - f_{Y_{-\infty};0}^\theta| \leq b_1^{m+1} x + \sum_{j=m+1}^{\infty} b_1^j (a_0 + a_1 Y_{-j}^2) \xrightarrow{\mathbb{P}_* \text{-a.s.}} 0,$$

showing **H 3.23**-(a). Note similarly that

$$\begin{aligned} & \left| \log q^\theta(f_{Y_{1:t-1}}^\theta(x), Y_t) - \log q^\theta(f_{Y_{-\infty:t-1}}^\theta, Y_t) \right| \\ &= \frac{\nu+1}{2} \left| \ln \left( 1 + \frac{Y_t^2}{\nu(f_{Y_{1:t-1}}^\theta(x))^2} \right) - \ln \left( 1 + \frac{Y_t^2}{\nu(f_{Y_{-\infty:t-1}}^\theta)^2} \right) \right|. \end{aligned}$$

Since for any  $a > 0$ ,  $b > 0$  and  $z > 0$ ,  $|\log(1+z^2/a^2) - \log(1+z^2/b^2)| \leq 2|a-b|/(a \wedge b)$  and  $f_{Y_{1:t-1}}^\theta(x) > \alpha_0$ ,  $f_{Y_{-\infty:t-1}}^\theta > \alpha_0$ , there exists a constant  $K > 0$  such that

$$\begin{aligned} & \sup_{\theta \in \Theta} \left| \log q^\theta(f_{Y_{1:t-1}}^\theta(x), Y_t) - \log q^\theta(f_{Y_{-\infty:t-1}}^\theta, Y_t) \right| \\ & \leq K \left( b_1^{m+1} x + \sum_{j=m+1}^{\infty} b_1^j (a_0 + a_1 Y_{-j}^2) \right) \xrightarrow{\mathbb{P}_* \text{-a.s.}} 0, \end{aligned}$$

showing **H 3.23**-(b). The proof of **H 3.23**-(c) is along the same lines. Theorem 3.25 shows the consistency of the estimator.

Of course, in well-specified models where the observations process is associated to a parameter  $\theta_* \in \Theta$ , Theorem 3.25 allows us to obtain the convergence of the MLE to  $\theta_*$  on the condition that  $\Theta_*$  is reduced to the singleton  $\{\theta_*\}$ .

An important subclass of misspecified models corresponds to the case where the kernel density  $q^\theta = q$  does not depend on  $\theta$  and where the observation process is assumed to follow the recursions:

$$\begin{aligned} \mathbb{P}[Y_t \in A | \mathcal{F}_{t-1}] &= Q^*(X_{t-1}, A) = \int_A q^*(X_{t-1}, y) \mu(dy), \quad \text{for any } A \in \mathcal{Y}, \\ X_t &= f_{Y_t}^{\theta_*}(X_{t-1}), \quad t \in \mathbb{Z}. \end{aligned}$$

The MLE is based on the likelihood functions  $y_{1:n} \mapsto L_{n,x}^\theta(y_{1:n})$  defined (3.38) whose expression includes a kernel density  $q \neq q^*$  and a family of functions  $\{f^\theta : \theta \in \Theta\}$ . Since  $\theta_*$  is assumed to be in  $\Theta^\circ$ , the interior of  $\Theta$  but  $q^* \neq q$ , this situation falls into the misspecified models framework. In such cases, Maximum Likelihood Estimators  $\{\hat{\theta}_{n,x}\}$  defined in (3.37) are called Quasi Maximum Likelihood estimators (QMLE) and  $\theta_*$  is not the true value of the parameter in the sense that the distribution of the observation process is not characterized by  $\theta_*$  only. Nevertheless, perhaps surprisingly, it can be shown that, under some additional assumptions, the QMLE  $\{\hat{\theta}_{n,x}\}$  are consistent and asymptotically normal with respect to the parameter  $\theta_*$ . For simplicity, we assume here that  $X = \mathbb{R}$ .

The main assumption that links  $q^*$  and  $q$  is the following:

**H 3.27** For all  $x^* \in \mathbb{R}$ ,

$$\arg \max_{x \in \mathbb{R}} \int Q^*(x^*, dy) \ln q(x, y) = \{x^*\}. \quad (3.45)$$



**Theorem 3.28.** Assume **H 3.21**, **H 3.22**, **H 3.23** and **H 3.27**. Moreover, assume that  $f_{Y_{-\infty};0}^\theta = f_{Y_{-\infty};0}^{\theta_\star}$ ,  $\mathbb{P}$  – a.s. implies that  $\theta = \theta_\star$ . Then, for all  $x \in \mathbb{X}$ ,

$$\lim_{n \rightarrow \infty} d(\hat{\theta}_{n,x}, \theta_\star) = 0, \quad \mathbb{P} - \text{a.s.}$$

PROOF. See Exercise 3.52. ■

**Example 3.29 (Normal innovations)** Assume that the observations  $\{Y_t, t \in \mathbb{Z}\}$  is a strict-sense stationary ergodic process associated to

$$\begin{aligned} \mathbb{P}[Y_t \in A | \mathcal{F}_{t-1}] &= Q^\star(\sigma_{t-1}^2, A), \quad \text{for any } A \in \mathcal{Y}, \\ \sigma_t^2 &= f_{Y_t}^{\theta_\star}(\sigma_{t-1}^2), \quad t \in \mathbb{Z}. \end{aligned} \tag{3.46}$$

We fit to the observations an observation-driven model with normal innovations,

$$\begin{aligned} Y_t &= \sigma_t \varepsilon_t \\ \sigma_t^2 &= f_{Y_t}^\theta(\sigma_{t-1}^2), \quad (t, \theta) \in \mathbb{N} \times \Theta, \end{aligned} \tag{3.47}$$

where  $\{\varepsilon_t, t \in \mathbb{N}\} \sim_{\text{iid}} \mathcal{N}(0, 1)$ . This model typically covers the GARCH(1,1) model with normal innovations. Now, the function

$$x \mapsto \int Q^\star(x^\star, dy) \ln q(x, y) = \int Q^\star(x^\star, dy) \left( -\frac{y^2}{2x} - \frac{1}{2} \ln(2\pi x) \right) = \left( -\frac{\int Q^\star(x^\star, dy) y^2}{2x} - \frac{1}{2} \ln(2\pi x) \right),$$

is maximized at  $x = \int Q^\star(x^\star, dy) y^2$  by straightforward algebra. Thus, **H 3.27** implies

$$\int Q^\star(x^\star, dy) y^2 = x^\star.$$

Plugging this equality into (3.46), we obtain that the observations  $\{Y_t, t \in \mathbb{Z}\}$  is a strict-sense stationary ergodic process associated to

$$\begin{aligned} Y_t &\sim \sigma_t \varepsilon_t^\star \\ \sigma_t^2 &= f_{Y_t}^{\theta_\star}(\sigma_{t-1}^2), \quad t \in \mathbb{Z}, \end{aligned}$$

where  $\{\varepsilon_t^\star, t \in \mathbb{Z}\}$  is an i.i.d. sequence of random variables with potentially any unknown distribution provided that  $\mathbb{E}[(\varepsilon_t^\star)^2] = 1$ .

**Example 3.30 (Exponential families)** Assume that the observations  $\{Y_t, t \in \mathbb{Z}\}$  are a strict-sense stationary ergodic process associated to

$$\begin{aligned} \mathbb{P}[Y_t \in A | \mathcal{F}_{t-1}] &= Q^\star(X_{t-1}, A) = \int_A q^\star(X_{t-1}, y) \mu(dy), \quad \text{for any } A \in \mathcal{Y}, \\ X_t &= f_{Y_t}^{\theta_\star}(X_{t-1}), \quad t \in \mathbb{Z}. \end{aligned}$$

We fit to the observations the following observation-driven model

$$\begin{aligned} \mathbb{P}[Y_t \in A | \mathcal{F}_{t-1}] &= Q(X_{t-1}, A), \quad \text{for any } A \in \mathcal{Y}, \\ X_t &= f_{Y_t}^\theta(X_{t-1}), \quad (t, \theta) \in \mathbb{Z} \times \Theta. \end{aligned}$$

where  $Q(x, \cdot)$  is assumed to belong to the class of exponential family distributions. More precisely, we assume that for all  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ ,  $q(x, y) = \exp(xy - A(x))h(y)$  for some twice differentiable function  $A : \mathbb{X} \rightarrow \mathbb{R}$  and some measurable function  $h : \mathbb{Y} \rightarrow \mathbb{R}^+$ . Using that

$$\int Q(x, dy) \frac{\partial^2 \ln q(x, y)}{\partial x^2} \leq 0,$$

it can be readily checked that  $A'' \geq 0$  so that  $A$  is concave. Thus, the function

$$\begin{aligned} x \mapsto \int Q^*(x^*, dy) \ln q(x, y) &= \int Q^*(x^*, dy) (xy - A(x) + \ln h(y)) \\ &= x \int Q^*(x^*, dy) y - A(x) + \int Q^*(x^*, dy) \ln h(y), \end{aligned}$$

is convex. The maximum of this function can thus be obtained by cancelling the derivatives with respect to  $x$ , which yields  $\int Q^*(x^*, dy) y - A'(x) = 0$ . Then, **H 3.27** implies that

$$\int Q^*(x^*, dy) y = A'(x^*). \quad (3.48)$$

For example in the log-linear Poisson autoregression model (see Example 2.48),

$$q(x, y) = \exp(xy - e^x)/y!$$

so that  $A(x) = \exp(x)$ . Thus, (3.48) yields that  $\int Q^*(x^*, dy) y = \exp(x^*)$ .

**Remark 3.31** In well-specified models,  $Q^* = Q$ . Since the Kullback-Leibler divergence is nonnegative, we obtain

$$\int Q^*(x^*, dy) \ln q^*(x, y) \leq \int Q^*(x^*, dy) \ln q^*(x^*, y),$$

and, provided that  $x \mapsto Q(x, \cdot)$  is one-to-one, the equality holds if and only if  $x = x^*$ . Thus, (3.45) in **H 3.27** is satisfied.

If for all  $y \in \mathbb{Y}$ ,  $x \mapsto q(x, y)$  is twice differentiable, we may define for all vector  $\mathbf{u}$  in  $\mathbb{R}^d$ , all matrix  $\Gamma$  of size  $d \times d$  with real entries and all  $(x, y) \in \mathbb{R} \times \mathbb{Y}$ ,

$$\varphi(\mathbf{u}, x, y) := \mathbf{u} \frac{\partial \ln q}{\partial x}(x, y), \quad (3.49)$$

$$\psi(\Gamma, \mathbf{u}, x, y) := \Gamma \frac{\partial \ln q}{\partial x}(x, y) + \mathbf{u} \mathbf{u}' \frac{\partial^2 \ln q}{\partial x^2}(x, y). \quad (3.50)$$

These functions appear naturally when differentiating  $\theta \mapsto \ln q(f(\theta), y)$  where  $\theta \mapsto f(\theta)$  is a twice differentiable function. More precisely, we have in such a case by straightforward algebra,

$$\begin{aligned} \nabla \ln q(f(\theta), y) &= \varphi(\nabla f(\theta), f(\theta), y), \\ \nabla^2 \ln q(f(\theta), y) &= \psi(\nabla^2 f(\theta), \nabla f(\theta), f(\theta), y). \end{aligned}$$

**H 3.32** For all  $y \in \mathbb{Y}$ , the function  $x \mapsto q(x, y)$  is twice continuously differentiable. Moreover, there exist  $\rho > 0$  and a family of  $\mathbb{P}$ -a.s. finite random variables

$$\left\{ f_{Y-\infty:t}^\theta : (\theta, t) \in \Theta \times \mathbb{Z} \right\}$$

such that  $\theta \mapsto f_{Y-\infty:0}^\theta$  is,  $\mathbb{P}$ -a.s., twice continuously differentiable on some ball  $B(\theta_*, \rho)$  and for all  $x \in \mathbb{X}$ ,

(a)  $\mathbb{P}$ -a.s.,

$$\lim_{t \rightarrow \infty} \|\varphi(\nabla f_{Y_{1:t-1}}^{\theta_*}(x), f_{Y_{1:t-1}}^{\theta_*}(x), Y_t) - \varphi(\nabla f_{Y_{-\infty:t-1}}^{\theta_*}, f_{Y_{-\infty:t-1}}^{\theta_*}, Y_t)\| = 0,$$

where  $\|\cdot\|$  is any norm on  $\mathbb{R}^d$ ,

(b)  $\mathbb{P}$ -a.s. ,

$$\lim_{t \rightarrow \infty} \sup_{\theta \in B(\theta_*, \rho)} \|\psi(\nabla^2 f_{Y_{1:t-1}}^\theta(x), \nabla f_{Y_{1:t-1}}^\theta(x), f_{Y_{1:t-1}}^\theta(x), Y_t) - \psi(\nabla^2 f_{Y_{-\infty:t-1}}^\theta, \nabla f_{Y_{-\infty:t-1}}^\theta, f_{Y_{-\infty:t-1}}^\theta, Y_t)\| = 0,$$

where by abuse of notation, we use again  $\|\cdot\|$  to denote any norm on the set of  $d \times d$ -matrices with real entries,

(c)

$$\mathbb{E} \left[ \|\varphi(\nabla f_{Y_{-\infty;0}}^{\theta_*}, f_{Y_{-\infty;0}}^{\theta_*}, Y_1)\|^2 \right] < \infty, \quad (3.51)$$

$$\mathbb{E} \left[ \sup_{\theta \in B(\theta_*, \rho)} \|\psi(\nabla^2 f_{Y_{-\infty;0}}^\theta, \nabla f_{Y_{-\infty;0}}^\theta, f_{Y_{-\infty;0}}^\theta, Y_1)\| \right] < \infty. \quad (3.52)$$

Moreover, the matrix  $\mathcal{J}(\theta_*)$  defined by

$$\mathcal{J}(\theta_*) := \mathbb{E} \left[ (\nabla f_{Y_{-\infty;0}}^{\theta_*})(\nabla f_{Y_{-\infty;0}}^{\theta_*})' \left( \frac{\partial^2 \ln q}{\partial x^2}(f_{Y_{-\infty;0}}^{\theta_*}, y) \right) \right], \quad (3.53)$$

is nonsingular.

**Theorem 3.33.** Assume **H** 3.21, **H** 3.27 and **H** 3.32. Assume in addition that  $\hat{\theta}_{n,x} \rightarrow_{\mathbb{P}} \theta_*$ . Then,

$$\sqrt{n}(\hat{\theta}_{n,x} - \theta_*) \xrightarrow{\mathcal{L}_{\mathbb{P}}} N(0, \mathcal{J}(\theta_*)^{-1} \mathcal{I}(\theta_*) \mathcal{J}(\theta_*)^{-1}),$$

where  $\mathcal{J}(\theta_*)$  is given by (3.53) and  $\mathcal{I}(\theta_*)$  is defined by

$$\mathcal{I}(\theta_*) := \mathbb{E} \left[ (\nabla f_{Y_{-\infty;0}}^{\theta_*})(\nabla f_{Y_{-\infty;0}}^{\theta_*})' \left( \frac{\partial \ln q}{\partial x}(f_{Y_{-\infty;0}}^{\theta_*}, y) \right)^2 \right].$$

We preface the proof by two Lemmas, whose proofs are given as exercises.

**Lemma 3.34** Assume **H** 3.27 and **H** 3.32. Then,

$$n^{-1/2} \sum_{t=1}^n \nabla \ln q(f_{Y_{-\infty:t-1}}^{\theta_*}, Y_t) \xrightarrow{\mathcal{L}_{\mathbb{P}}} N(0, \mathcal{I}(\theta_*)).$$

PROOF. See Exercise 3.53. ■

**Lemma 3.35** Assume **H** 3.27 and **H** 3.32. Let  $\{\theta_n, n \in \mathbb{N}\}$  be a sequence of random vectors such that  $\theta_n \rightarrow_{\mathbb{P}} \theta_*$ . Then, for all  $i, j \in \{1, \dots, d\}$ ,

$$n^{-1} \sum_{t=1}^n \frac{\partial^2 \ln q(f_{Y_{-\infty:t-1}}^{\theta_n}, Y_t)}{\partial \theta_i \partial \theta_j} \xrightarrow{\mathbb{P}\text{-prob}} \mathbb{E} \left[ \frac{\partial^2 \ln q(f_{Y_{-\infty:t-1}}^{\theta_*}, Y_t)}{\partial \theta_i \partial \theta_j} \right].$$

PROOF. See Exercise 3.54. ■

PROOF. [of Theorem 3.33] Since the MLE  $\hat{\theta}_{n,x}$  cancels the derivatives of the log-likelihood, a Taylor expansion at  $\theta = \theta_*$  with an integral form of the remainder yields

$$n^{-1/2} \sum_{t=1}^n \nabla \ln q(f_{Y_{1:t-1}}^{\hat{\theta}_{n,x}}(x), Y_t) = 0 = n^{-1/2} \sum_{t=1}^n \nabla \ln q(f_{Y_{1:t-1}}^{\theta_*}(x), Y_t) + n^{-1} \sum_{t=1}^n \left( \int_0^1 \nabla^2 \ln q(f_{Y_{1:t-1}}^{\theta_{n,x,s}}(x), Y_t) ds \right) \sqrt{n}(\hat{\theta}_n - \theta_*), \quad (3.54)$$

where  $\theta_{n,x,s} = s\hat{\theta}_{n,x} + (1-s)\theta_*$ . The proof of Theorem 3.33 then follows from (3.54) and the Slutsky Lemma, provided we show that for all  $\theta_n \rightarrow_{\mathbb{P}} \theta_*$ ,

$$n^{-1/2} \sum_{t=1}^n \nabla \ln q(f_{Y_{1:t-1}}^{\theta_*}(x), Y_t) \xrightarrow{\mathcal{L}_{\mathbb{P}}} N(0, \mathcal{J}(\theta_*)), \quad (3.55)$$

$$n^{-1} \sum_{t=1}^n \frac{\partial^2 \ln q(f_{Y_{1:t-1}}^{\theta_n}(x), Y_t)}{\partial \theta_i \partial \theta_j} \xrightarrow{\mathbb{P}\text{-prob}} \mathbb{E} \left[ \frac{\partial^2 \ln q(f_{Y_{-\infty:0}}^{\theta_*}, Y_1)}{\partial \theta_i \partial \theta_j} \right], \quad (3.56)$$

$$\mathcal{J}(\theta_*) = \mathbb{E} \left[ \psi(\nabla^2 f_{Y_{-\infty:0}}^{\theta_*}, \nabla f_{Y_{-\infty:0}}^{\theta_*}, f_{Y_{-\infty:0}}^{\theta_*}, Y_1) \right]. \quad (3.57)$$

By noting that

$$\begin{aligned} \nabla \ln q(f_{Y_{1:t-1}}^{\theta_*}(x), Y_t) &= \varphi(\nabla f_{Y_{1:t-1}}^{\theta_*}(x), f_{Y_{1:t-1}}^{\theta_*}(x), Y_t), \\ \nabla^2 \ln q(f_{Y_{1:t-1}}^{\theta_*}(x), Y_t) &= \psi(\nabla^2 f_{Y_{1:t-1}}^{\theta_*}(x), \nabla f_{Y_{1:t-1}}^{\theta_*}(x), f_{Y_{1:t-1}}^{\theta_*}(x), Y_t), \end{aligned}$$

we obtain, according to **H** 3.32-(a)-(b), that showing (3.55)-(3.56) is equivalent to showing that for all  $\theta_n \rightarrow_{\mathbb{P}} \theta_*$ ,

$$n^{-1/2} \sum_{t=1}^n \nabla \ln q(f_{Y_{-\infty:t-1}}^{\theta_*}, Y_t) \xrightarrow{\mathcal{L}_{\mathbb{P}}} N(0, \mathcal{J}(\theta_*)), \quad (3.58)$$

$$n^{-1} \sum_{t=1}^n \frac{\partial^2 \ln q(f_{Y_{-\infty:t-1}}^{\theta_n}, Y_t)}{\partial \theta_i \partial \theta_j} \xrightarrow{\mathbb{P}\text{-prob}} \mathbb{E} \left[ \frac{\partial^2 \ln q(f_{Y_{-\infty:0}}^{\theta_*}, Y_1)}{\partial \theta_i \partial \theta_j} \right]. \quad (3.59)$$

But this follows from Lemma 3.34 and Lemma 3.35. It remains to show (3.57). Since under **H** 3.27,

$$\mathbb{E} \left[ \nabla^2 f_{Y_{-\infty:0}}^{\theta_*} \frac{\partial \ln q}{\partial x}(\nabla f_{Y_{-\infty:0}}^{\theta_*}, Y_1) \right] = 0,$$

we have, using (3.50),

$$\mathcal{J}(\theta_*) = \mathbb{E} \left[ \psi(\nabla^2 f_{Y_{-\infty:0}}^{\theta_*}, \nabla f_{Y_{-\infty:0}}^{\theta_*}, f_{Y_{-\infty:0}}^{\theta_*}, Y_1) \right].$$

The proof is completed. ■

### 3.4 Bayesian inference

In addition to classical or likelihood inference, there is another approach termed Bayesian inference, wherein prior belief is combined with data to obtain posterior distributions on which statistical inference is based. Although we focus primarily on likelihood based inference in this text, it is also worthwhile to consider Bayesian inference for some problems. Except for some simple cases, Bayesian inference can be computationally intensive and may rely on computational techniques.

The basic idea in Bayesian analysis is that a parameter vector, say  $\theta \in \Theta$ , is unknown to a researcher, so a *prior* distribution,  $\pi(\theta)$ , is put on the parameter vector. The researcher also proposes a model or likelihood,  $p(x | \theta)$ , that describes how the data  $X$  depend on the parameter vector. Inference about  $\theta$  is then based on the *posterior* distribution, which is obtained via Bayes's theorem,

$$\pi(\theta | X = x) \propto \pi(\theta) p(x | \theta).$$

In some simple cases, the prior and the likelihood are *conjugate* distributions that may be combined easily. For example, in  $n$  fixed repeated (i.i.d.) Bernoulli experiments with probability of success  $\theta$ , a *Beta-Binomial* conjugate pair is taken. In this case the prior is  $\text{Beta}(a, b)$ :  $\pi(\theta) \propto \theta^a (1 - \theta)^b$ ; the values  $a, b > -1$  are called hyperparameters. The likelihood in this example is  $\text{Binomial}(n, \theta)$ :  $p(x | \theta) \propto \theta^x (1 - \theta)^{n-x}$ , from which we easily deduce that the posterior is also Beta,  $\pi(\theta | X = x) \propto \theta^{x+a} (1 - \theta)^{n+b-x}$ , and from which inference may easily be achieved. As previously mentioned, in more complex experiments, the posterior distribution is often difficult to obtain by direct calculation, so MCMC techniques are employed. The main idea is that we may not be able to explicitly display the posterior, but we may be able to simulate from the posterior.

In the remaining part of this section, we give three examples of the Bayesian analysis of time series models. The first example is a simple linear  $\text{AR}(p)$  model. Then, we discuss two nonlinear examples, a threshold model and a GARCH model.

**Example 3.36 (AR(p) Model)** We now discuss the Bayesian approach for fitting a normal autoregressive model of order  $p$ . We write the model as  $Y_t = \sum_{i=1}^p \phi_i Y_{t-i} + \sigma Z_t$  where  $\{Z_t, t \in \mathbb{N}\} \sim_{\text{iid}} N(0, 1)$ ; replace  $Y_s$  by  $Y_s - \mu$  if the mean is not zero. Define

$$\mathbf{Y} = \begin{pmatrix} Y_{p+1} \\ Y_{p+2} \\ \vdots \\ Y_n \end{pmatrix}, \quad \boldsymbol{\phi} = \begin{pmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_p \end{pmatrix}, \quad \mathbf{Z} = \begin{pmatrix} Z_{p+1} \\ Z_{p+2} \\ \vdots \\ Z_n \end{pmatrix},$$

$$X = \begin{pmatrix} Y_p & Y_{p-1} & \dots & Y_1 \\ Y_{p+1} & Y_p & \dots & Y_2 \\ \vdots & \vdots & & \vdots \\ Y_n & Y_{n-1} & \dots & Y_{n-p+1} \end{pmatrix}.$$

The model can then be written as

$$\mathbf{Y} = X\boldsymbol{\phi} + \sigma\mathbf{Z},$$

which implies that  $\mathbf{Y} \sim N(X\boldsymbol{\phi}, \sigma^2\mathbf{I})$ , with  $X$  being  $(n-p) \times p$  matrix; for simplicity,  $X$  is assumed to be of full column rank  $p$ . It is easier to reparametrize the model in terms of a precision parameter,  $\tau = 1/\sigma^2$ . Then the likelihood function can be written as, up to a constant,

$$L_{\boldsymbol{\phi}, \tau} \mathbf{Y} = p(\mathbf{Y} | \boldsymbol{\phi}, \tau) \propto \tau^{(n-p)/2} \exp(-\tau(\mathbf{Y} - X\boldsymbol{\phi})'(\mathbf{Y} - X\boldsymbol{\phi})/2). \quad (3.60)$$

In a Bayesian framework, there are two ways for considering the prior settings of the ARMA parameters. One is to constrain the prior to be non-zero only over the region defined by the stationarity conditions, such as in ? and ?. The second method is to simply ignore the stationary assumptions and proceed to use, e.g., the normal-gamma family as a conjugate prior distribution, such as in ?. In this example, we apply the latter method, i.e., no restriction for stationarity or invertibility is made in the prior distribution. We use the following specification

$$\begin{cases} \mathbf{Y} | (\boldsymbol{\phi}, \tau) \sim N(X\boldsymbol{\phi}, \mathbf{I}/\tau) \\ \boldsymbol{\phi} | \tau \sim N_p(\boldsymbol{\phi}_0, V_0/\tau) \\ \tau \sim \Gamma(a, b) \end{cases}$$

where  $\boldsymbol{\phi}_0$  and  $V_0$  are the prior mean and covariance matrix, and  $a$  and  $b$  are the shape and scale of the prior of the precision parameter. The matrix  $V_0$  is assumed to be invertible (which is generally chosen to be  $V_0 = \gamma_0^2 \mathbf{I}$ , and we can take  $\gamma_0^2$  large if the prior is meant to be noninformative). After straightforward simplification, the posterior joint density of  $(\boldsymbol{\phi}, \tau)$  becomes

$$\begin{aligned} \pi(\boldsymbol{\phi}, \tau | \mathbf{Y}) &\propto \tau^{\tilde{a}-1} \exp(-\tau(b + (1/2)\boldsymbol{\phi}_0' V_0^{-1} \boldsymbol{\phi}_0 + (1/2)\|\mathbf{y} - X\boldsymbol{\phi}\|^2)) \\ &\propto \tau^{\tilde{a}-1} \exp(-\tau\tilde{b}) \exp(-(\tau/2) \cdot (\boldsymbol{\phi} - \bar{\boldsymbol{\phi}})'(X'X + V_0^{-1})(\boldsymbol{\phi} - \bar{\boldsymbol{\phi}})) \end{aligned}$$

where

$$\begin{aligned}\tilde{a} &= \frac{n}{2} + a \\ \tilde{b} &= b + \frac{\mathbf{Y}'\mathbf{Y} - (\mathbf{X}'\mathbf{Y} + V_0^{-1}\boldsymbol{\phi}_0)'(\mathbf{X}'\mathbf{X} + V_0^{-1})^{-1}(\mathbf{X}'\mathbf{Y} + V_0^{-1}\boldsymbol{\phi}_0)}{2} \\ \bar{\boldsymbol{\phi}} &= (\mathbf{X}'\mathbf{X} + V_0^{-1})^{-1}(\mathbf{X}'\mathbf{Y} + V_0^{-1}\boldsymbol{\phi}_0).\end{aligned}$$

From the joint posterior density  $f(\boldsymbol{\phi}, \boldsymbol{\tau} \mid \mathbf{Y})$ , it can be seen that the conditional posterior density of  $(\boldsymbol{\phi} \mid \boldsymbol{\tau}, \mathbf{Y})$  is multivariate normal with mean  $\bar{\boldsymbol{\phi}}$  and covariance matrix  $\boldsymbol{\tau}^{-1}(\mathbf{X}'\mathbf{X} + V_0^{-1})^{-1}$ , i.e.,

$$\boldsymbol{\phi} \mid (\boldsymbol{\tau}, \mathbf{Y}) \sim N_p(\bar{\boldsymbol{\phi}}, \boldsymbol{\tau}^{-1}(\mathbf{X}'\mathbf{X} + V_0^{-1})^{-1}).$$

The conditional posterior density of  $(\boldsymbol{\tau} \mid \boldsymbol{\phi}, \mathbf{Y})$  is  $\Gamma(\tilde{a}, \tilde{b}(\boldsymbol{\phi}))$ ,

$$\tilde{b}(\boldsymbol{\phi}) := b + (1/2)\boldsymbol{\phi}_0'V_0^{-1}\boldsymbol{\phi}_0 + (1/2)\|\mathbf{Y} - \mathbf{X}\boldsymbol{\phi}\|^2.$$

As an example, we use our R script `arp.mcmc` to analyze a simulated AR(2) model with  $\phi_1 = 1$ ,  $\phi_2 = -0.9$  and  $\sigma^2 = 1$ .

```
> x = arima.sim(list(order=c(2,0,0), ar=c(1,-.9)), n=1000) # generate an ar2
> arp.mcmc(x, porder=2, n.iter=1000, n.warmup=100) # sample from posterior
```

phi.Series 1		phi.Series 2		tauinv	
Min.	:0.9831	Min.	:-0.9598	Min.	:0.8068
1st Qu.	:1.0142	1st Qu.	:-0.9292	1st Qu.	:0.9117
Median	:1.0228	Median	:-0.9219	Median	:0.9419
Mean	:1.0226	Mean	:-0.9217	Mean	:0.9430
3rd Qu.	:1.0311	3rd Qu.	:-0.9138	3rd Qu.	:0.9739
Max.	:1.0672	Max.	:-0.8744	Max.	:1.0886

A summary as well as trace plots are produced by the script; see Figure 3.1.

**Example 3.37 (SETAR model for the lynx series)** A two-regime SETAR model is considered, and denoted by SETAR (2;  $p_1, p_2$ ):

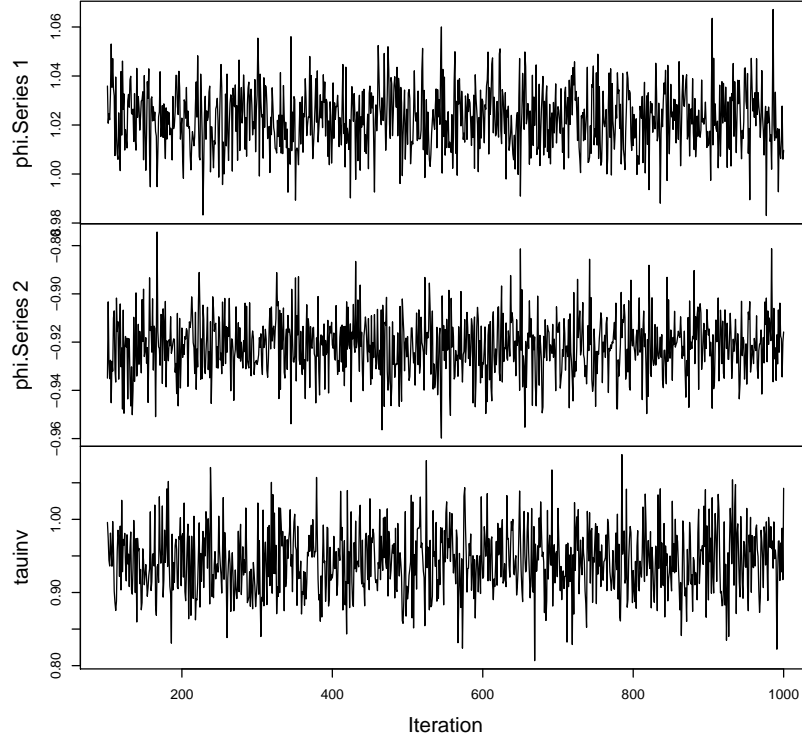
$$Y_t = \begin{cases} \phi_0^{(1)} + \sum_{i=1}^{p_1} \phi_i^{(1)} Y_{t-i} + Z_t^{(1)} & Y_{t-d} \leq r \\ \phi_0^{(2)} + \sum_{i=1}^{p_2} \phi_i^{(2)} Y_{t-i} + Z_t^{(2)} & Y_{t-d} > r \end{cases} \quad (3.61)$$

where  $Z_t^{(k)} = \sigma_k Z_t$ ,  $k = 1, 2$ , and  $\{Z_t, t \in \mathbb{N}\} \sim_{\text{iid}} N(0, 1)$ . This is a piecewise AR model, where regime switching is driven by lagged values of the series and a threshold value  $r$ . The unknown parameters in the model are  $\boldsymbol{\phi}_k = (\phi_0^{(k)}, \phi_1^{(k)}, \dots, \phi_{p_k}^{(k)})'$ ,  $k = 1, 2$ , the threshold value  $r$ , the delay lag  $d$  and the regime error variances are  $\sigma_k^2$ ,  $k = 1, 2$ .

Likelihood approaches for estimating threshold models usually involve setting the threshold  $r$  and the delay  $d$  (e.g.,  $r = 0$  and  $d = 1$ ), or choosing  $r, d$  via some information criterion. The likelihood inference for the parameters is then carried out conditionally upon these choices, but ignore the associated uncertainty in those parameters. The uncertainty in estimating  $r$  and  $d$  is explicitly captured via a Bayesian approach.

For SETAR models, the stationary and invertible conditions in the literature are rather involved. As in Example 3.36, we do not take these constraints into account in the Bayesian analysis. We use the following specification for the priors. The AR parameters are assumed to be normal with zero mean and covariance  $V_k$ ,  $k = 1, 2$ , i.e.,  $\boldsymbol{\phi}_k \sim N(0, V_k)$ , where  $V_k$  is the prior covariance matrix. We may for example set  $V_k = \gamma_v^2 \mathbf{I}_{g_k}$ , where  $\gamma_v^2$  is the prior variance (generally taken to be large). An inverse-gamma prior is assumed for  $\sigma_k^2$ , i.e.,  $1/\sigma_k^2 \sim \Gamma(a_k, b_k)$ , for  $k = 1, 2$ , where the hyper-parameters  $a_k$  and  $b_k$  are known constants. The prior for the threshold value  $r$  is uniformly distributed between  $[r_{\min}, r_{\max}]$ , where  $r_{\min}$  and  $r_{\max}$  are generally specified as percentiles of the observations (e.g., 10 and 90 percentiles), for identification purposes in each regime. A discrete uniform prior on  $\{1, 2, \dots, d_{\max}\}$  is employed for the delay  $d$ .

Let  $\mathbf{Y} = (Y_{p+1}, \dots, Y_n)$  be the observed values of the TAR model, where  $p = \max(p_1, p_2)$ . the conditional likelihood function of the model is:



**Fig. 3.1** Output of `arp.mcmc` for Example 3.36. Each plot shows the trace, after burnin, of the draws for each component.

$$L(\mathbf{y} \mid \boldsymbol{\phi}^{(1)}, \boldsymbol{\phi}^{(2)}, \sigma_1^2, \sigma_2^2, r, d) \propto \sigma_1^{-n_1} \sigma_2^{-n_2} \times \exp \left( \sum_{t=p+1}^n \sum_{k=1}^2 -\frac{1}{2\sigma_k^2} \left\{ Y_t - \phi_0^{(k)} - \sum_{i=1}^{p_k} \phi_j^{(k)} Y_{t-i} \right\}^2 I_t^{(k)} \right)$$

where  $n_1 := \sum_{t=p+1}^n \mathbb{1}\{Y_{t-d} \leq r\}$ ,  $n_2 = n - n_1$  and  $I_t^{(k)}$  is the indicator variable  $\mathbb{1}\{r_{k-1} \leq Y_{t-d} < r_k\}$  with  $r_0 = -\infty$ ,  $r_1 = r$  and  $r_2 = \infty$ . For notational convenience, let the full parameter vector be  $\boldsymbol{\theta} = \{\boldsymbol{\phi}^{(1)}, \boldsymbol{\phi}^{(2)}, \sigma_1^2, \sigma_2^2, r, d\}$ . Under the prior specification, the  $k$ -th regime AR parameter  $\boldsymbol{\phi}^{(k)} \mid (\mathbf{Y}, \boldsymbol{\theta} \setminus \{\boldsymbol{\phi}^{(k)}\})$  is Gaussian with mean  $\bar{\boldsymbol{\phi}}^{(k)}$  and covariance  $\bar{V}_k$  given by

$$\bar{\boldsymbol{\phi}}^{(k)} = \sigma_k^{-2} \bar{V}_k X_k' \mathbf{Y}_k$$

$$\bar{V}_k = (\sigma_k^{-2} X_k' X_k + V_k^{-1})^{-1},$$

where the design matrices  $X_1$  and  $X_2$  are given by

$$X_1 = \begin{bmatrix} 1 & Y_{\ell_1-1} & \cdots & \cdots & Y_{\ell_1-p_1} \\ 1 & Y_{\ell_2-1} & \cdots & \cdots & Y_{\ell_2-p_1} \\ \vdots & \vdots & & & \vdots \\ 1 & Y_{\ell_{n_1}-1} & \cdots & \cdots & Y_{\ell_{n_1}-p_1} \end{bmatrix}$$

and

$$X_2 = \begin{bmatrix} 1 & Y_{u_1-1} & \cdots & \cdots & Y_{u_1-p_2} \\ 1 & Y_{u_2-1} & \cdots & \cdots & Y_{u_2-p_2} \\ \vdots & \vdots & & & \vdots \\ 1 & Y_{u_{n_1}-1} & \cdots & \cdots & Y_{u_{n_1}-p_2} \end{bmatrix}$$

where  $\ell_i$  and  $u_j$  are the time indices of the  $i$ -th and  $j$ -th observations in the “lower” and “upper” regimes:  $\ell_i$  is the index for regime 1 and  $u_j$  is for regime 2. The vectors  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$  are set as  $\mathbf{Y}_1 = (Y_{\ell_1}, \dots, Y_{\ell_{n_1}})'$  and  $\mathbf{Y}_2 = (Y_{u_1}, \dots, Y_{u_{n_2}})'$ .

The posterior distribution of the inverse of the variance  $\sigma_k^{-2} \mid (\mathbf{Y}, \boldsymbol{\theta} \setminus \{\sigma_k^{-2}\})$  is Gamma with shape and scale

$$\tilde{a}_k = a_k + n_k/2, \quad \tilde{b}_k = b_k + (1/2)\|\mathbf{Y}_k - X_k \boldsymbol{\phi}_k\|^2.$$

The posterior distribution for the threshold is a bit more involved. By simple multiplication of the likelihood and prior for  $r$ , the posterior density of  $r \mid (\mathbf{Y}, \boldsymbol{\theta} \setminus \{r\})$  is given by (up to a normalization constant)

$$\sigma_1^{-n_1} \sigma_2^{-n_2} \exp\left(-\frac{1}{2} \sum_{k=1}^2 \frac{\|\mathbf{Y}_k - X_k \boldsymbol{\phi}_k\|^2}{\sigma_k^2}\right) \mathbb{1}_{\{r \in [r_{\min}, r_{\max}]\}}$$

where  $n_1, n_2, \mathbf{Y}_1, \mathbf{Y}_2, X_1, X_2$  depend on  $r$ . This posterior distribution is not a known or standard form and therefore we cannot use the elementary Gibbs sampler. To solve this problem, we use a hybrid of the ideas behind the Gibbs sampler and the Metropolis-Hastings algorithm. The algorithm merely replaces the step that samples from the full-posterior by a Metropolis-Hastings step, which leaves the full-posterior stationary. Now by simulating sufficiently many of these Metropolis-Hastings steps before proceeding to the next component, we will approximately mimic the elementary Gibbs sampler itself. However, there is often a far more efficient procedure that only updates the component once (or a small finite number of times) before moving on to other coordinates. Like the Gibbs sampler, this algorithm can be thought of as a special case of the Metropolis-Hastings algorithm. In this special case, we use a random walk Metropolis algorithm with a Gaussian proposal density with zero mean and variance  $\sigma^2$ , to sample this parameter. Finally, the posterior distribution of the delay  $d \mid (\mathbf{Y}, \boldsymbol{\theta} \setminus \{d\})$  is a multinomial trial with probabilities

$$\frac{L(\mathbf{Y} \mid \boldsymbol{\phi}^{(1)}, \boldsymbol{\phi}^{(2)}, \sigma_1^2, \sigma_2^2, r, j)}{\sum_{i=1}^{d_{\max}} L(\mathbf{Y} \mid \boldsymbol{\phi}^{(1)}, \boldsymbol{\phi}^{(2)}, \sigma_1^2, \sigma_2^2, r, i)}, \quad j = 1, 2, \dots, d_{\max}.$$

The sampling scheme amounts to updating parameters in turn, which leads to the algorithm implemented in the R package BAYSTAR. However, sampling the delay distribution this way is not optimal and may lead to slow mixing.

We analyzed the lynx data set presented in ???. We used BAYSTAR to perform a Bayesian inference of a two-regime threshold model, which was considered in ??. In particular, our goal was to fit the same model, an AR(2) for each regime, with the threshold delay being 2. In this example, we set the maximum delay,  $d_{\max}$ , equal to 2. In this case, the highest posterior probability for the delay lag was indeed 2. We also note that the variance of the proposal density,  $\sigma^2$ , which is called `step.thv`, must be specified. If the scaling parameter is too large, the acceptance ratio of the algorithm will be small. If it is too small, then the acceptance ratio is high but the algorithm mixes slowly (see Example 4.4 for a discussion of this issue). This can be monitored by observing the output produced by the package.

The R code and a partial output is given below. In addition to the output shown below, BAYSTAR exhibits various diagnostic plots, which we do not display.

```
> require(BAYSTAR)
> lynx.out = BAYSTAR(log10(lynx), lagp1=1:2, lagp2=1:2, Iteration=10000, Burnin=2000, d0=2,
  step.thv=.5)
      mean median   s.d. lower upper
phi1.0  0.5620 0.5701 0.1683 0.2039 0.8733
phi1.1  1.2614 1.2607 0.0724 1.1251 1.4022
phi1.2 -0.4143 -0.4174 0.0919 -0.5867 -0.2274
phi2.0  1.4567 1.5353 1.0021 -0.6318 3.2400
```



```

phi2.1      1.5694  1.5643  0.1317  1.3221  1.8346
phi2.2     -1.0623 -1.0759  0.2859 -1.5988 -0.4590
sigma^2 1    0.0363  0.0356  0.0067  0.0258  0.0513
sigma^2 2    0.0557  0.0535  0.0143  0.0349  0.0899
r           3.2057  3.2810  0.1818  2.5784  3.3848
acceptance rate of r = 23.76 %
The highest posterior prob. of lag is at : 2

```

We note that the result of the MCMC analysis is similar to the frequentist analysis presented in ??.

**Example 3.38 (Bayesian GARCH)** In this example, we consider a GARCH(1,1) model with Student- $t$  innovations for log-returns, i.e.,

$$Y_t = \sigma_t \left( \frac{v-2}{v} \right)^{1/2} \varepsilon_t \quad (3.62)$$

$$\sigma_t^2 = \alpha_0 + \alpha_1 Y_{t-1}^2 + \beta_1 \sigma_{t-1}^2 \quad (3.63)$$

where  $\{Z_t, t \in \mathbb{N}\}$  are i.i.d.  $t$ -distribution with  $v > 2$  degrees of freedom,  $\alpha_0 > 0$ ,  $\alpha_1, \beta_1 \geq 0$  and  $v > 2$ . To perform the Bayesian analysis of this model, it is easier to use data augmentation; see ? and ?. The idea is to represent the  $t$ -distribution as a scale mixture of Gaussians as suggested by ?. A random variable  $Z$  is a Gaussian scale mixture if it can be expressed as the product of a standard Gaussian  $G$  and an independent positive scalar random variable  $H^{1/2}$ :  $Z = H^{1/2}G$ . The variable  $H$  is the multiplier or the scale. If  $H$  has finite support, then  $Z$  is a finite mixture of Gaussians, whereas if  $H$  has a density (with respect to Lebesgue measure) on  $\mathbb{R}_+$ , then  $Z$  is a continuous mixture of Gaussian variables. Gaussian scale mixtures are symmetric, zero mean, and have leptokurtic marginal densities (tails heavier than those of a Gaussian distribution). Assume that the distribution of  $H$  is inverse gamma (denoted IG), with shape and scale parameters both equal to  $v/2$ ; then the distribution of  $H$  has the density

$$p_v(h) = \left( \frac{v}{2} \right)^{v/2} \frac{h^{-(v+2)/2}}{\Gamma(v/2)} \exp\left(-\frac{v}{2h}\right), h \geq 0.$$

If  $G \sim N(0, 1)$ , then the distribution of  $Z = H^{1/2}G$  has a density given by

$$f(z) = \frac{\left(\frac{v}{2}\right)^{v/2}}{\Gamma(v/2)\sqrt{2\pi}} \int_0^\infty h^{-(v+3)/2} \exp\left(-\frac{z^2 + v}{2h}\right) dh.$$

Using the result,

$$\int_0^\infty x^{-a/2} \exp\left(-\frac{b}{2x}\right) dx = \left(\frac{2}{b}\right)^{(a-2)/2} \Gamma\left(\frac{a-2}{2}\right),$$

which derives from a simple change variable in the definition of the complete gamma function, we obtain that

$$f(z) = \frac{\left(\frac{v}{2}\right)^{v/2} \Gamma((v+1)/2)}{\Gamma(v/2)\sqrt{2\pi}} \left(\frac{2}{z^2 + v}\right)^{(v+1)/2},$$

which is the density of a Student- $t$  distribution with  $v$  degrees of freedom. Using this representation of the  $t$ -distribution, we may rewrite (3.62) as follows,

$$\begin{aligned} Y_t &= \sigma_t \left( \frac{v-2}{v} \right)^{1/2} H_t^{1/2} G_t \\ G_t &\sim_{\text{iid}} N(0, 1) \\ H_t &\sim_{\text{iid}} \text{IG}\left(\frac{v}{2}, \frac{v}{2}\right) \\ \sigma_t^2 &= \alpha_0 + \alpha_1 Y_{t-1}^2 + \beta_1 \sigma_{t-1}^2. \end{aligned}$$

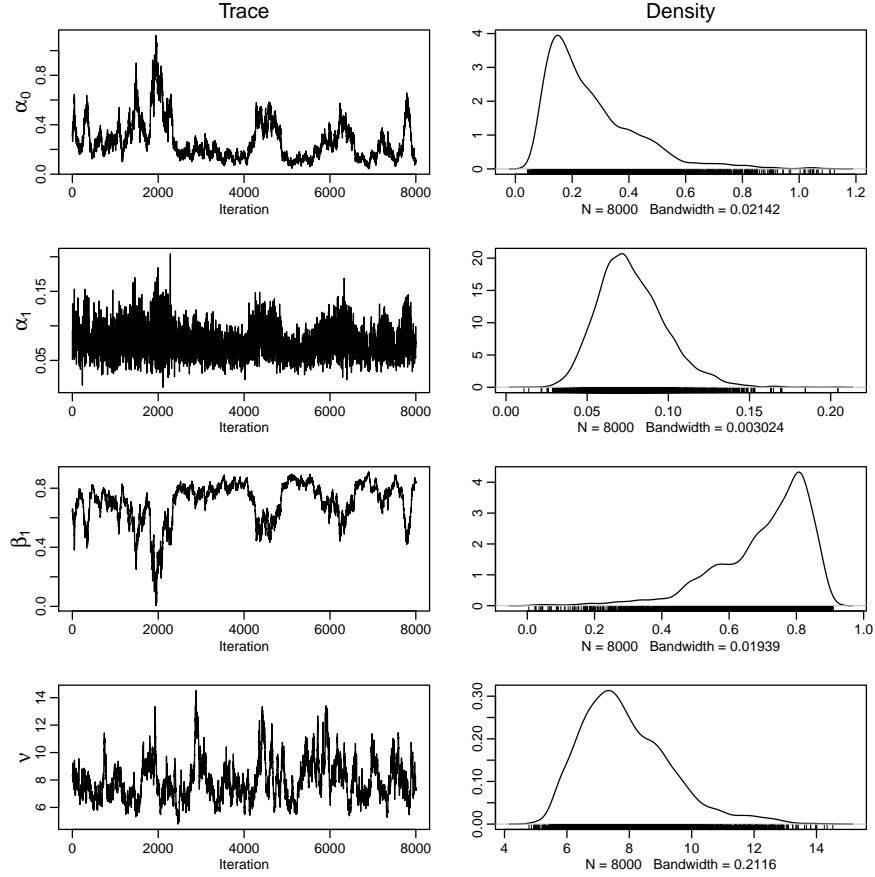


Fig. 3.2 Output of bayesGARCH for Example 3.38

Assume that we have  $n$  observations from this model. In order to write the likelihood function, we define the vectors,  $\mathbf{H} = (H_1, \dots, H_n)'$  and  $\boldsymbol{\alpha} = (\alpha_0, \alpha_1)'$ . The model parameters are collected into the vector  $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \beta, v)$ . Define the  $n \times n$  diagonal matrix

$$\Sigma = \Sigma(\boldsymbol{\theta}, \mathbf{H}) = \text{diag} \left( \left\{ H_t \frac{v-2}{v} \sigma_t^2(\boldsymbol{\alpha}, \beta_1) \right\}_{t=1}^n \right),$$

where  $\sigma_t^2(\boldsymbol{\alpha}, \beta_1) = \alpha_0 + \alpha_1 Y_{t-1}^2 + \beta_1 \sigma_{t-1}^2(\boldsymbol{\alpha}, \beta_1)$ , and  $\sigma_0^2(\boldsymbol{\alpha}, \beta_1) = 0$ . We can express the conditional likelihood

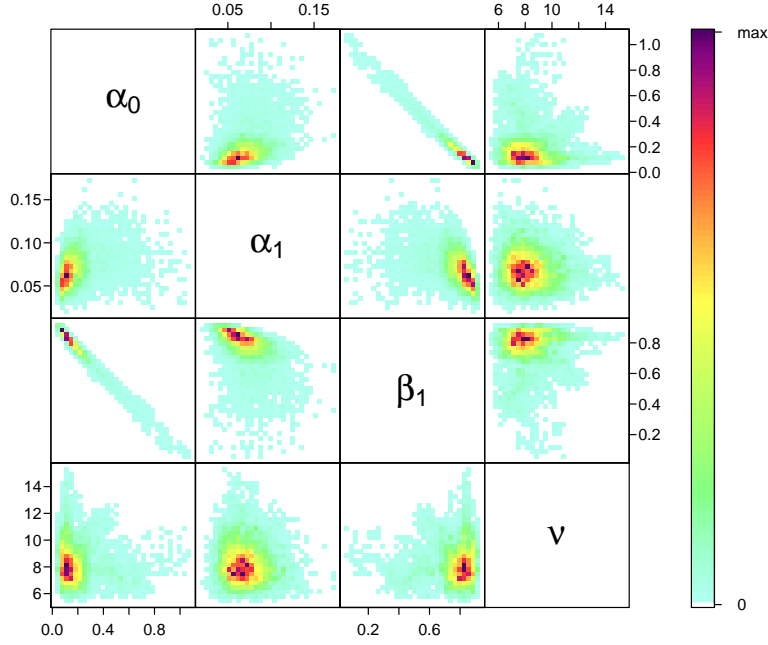
$$L(\mathbf{Y} \mid \boldsymbol{\theta}, \mathbf{H}) \propto (\det \Sigma)^{-1/2} \exp \left( -\frac{1}{2} \mathbf{Y}' \Sigma^{-1} \mathbf{Y} \right).$$

The prior distribution for the GARCH parameters  $\boldsymbol{\alpha}$  and  $\beta_1$  are chosen to be independent and distributed according to the truncated normal distribution (to ensure the positivity of the coefficients)

$$p(\boldsymbol{\alpha}) \propto (\det \Sigma_{\boldsymbol{\alpha}})^{-1/2} \exp \left( -\frac{1}{2} (\boldsymbol{\alpha} - \boldsymbol{\mu}_{\boldsymbol{\alpha}})' \Sigma_{\boldsymbol{\alpha}}^{-1} (\boldsymbol{\alpha} - \boldsymbol{\mu}_{\boldsymbol{\alpha}}) \right) \mathbb{I}_{\mathbb{R}_+^2}(\boldsymbol{\alpha})$$

$$p(\beta_1) \propto \sigma_{\beta}^{-1} \exp \left( -\frac{1}{2\sigma_{\beta}^2} (\beta_1 - \mu_{\beta})^2 \right) \mathbb{I}_{\mathbb{R}^+}(\beta_1),$$

where  $\mu_{\boldsymbol{\alpha}}$ ,  $\Sigma_{\boldsymbol{\alpha}}$ ,  $\mu_{\beta}$ ,  $\sigma_{\beta}$  are the hyperparameters. Only positivity constraints are taken into account in the prior specification; no stationarity conditions are imposed.



**Fig. 3.3** Pairwise sampling distributions for Example 3.38.

Since the components  $H_t$  are i.i.d. from the inverse gamma distribution, the conditional distribution of the vector  $\mathbf{H}$  given  $\mathbf{v}$  has a density given by

$$p(\mathbf{h} \mid \mathbf{v}) = \left(\frac{\mathbf{v}}{2}\right)^{\frac{n\mathbf{v}}{2}} \left[\Gamma\left(\frac{\mathbf{v}}{2}\right)\right]^{-n} \left(\prod_{t=1}^n h_t\right)^{-\frac{\mathbf{v}}{2}-1} \exp\left(-\frac{1}{2} \sum_{t=1}^n \frac{\mathbf{v}}{h_t}\right).$$

The prior distribution on the degrees of freedom  $\mathbf{v}$  is chosen to be a translated exponential with parameters  $\lambda > 0$  and  $\delta \geq 2$

$$p(\mathbf{v}) = \lambda \exp[-\lambda(\mathbf{v} - \delta)] \mathbb{I}\{\mathbf{v} \geq \delta\}.$$

For large values of  $\lambda$ , most of the mass of the prior is concentrated in the neighborhood of  $\delta$  and a constraint on the degrees of freedom can be imposed in this manner. A light tail prior may be imposed by taking both  $\delta$  and  $\lambda$  large. The joint prior distribution is then obtained by assuming that the priors are independent, i.e.,  $p(\tilde{\theta}, \tilde{h}) = p(\tilde{\alpha})p(\beta_1)p(h \mid \mathbf{v})p(\mathbf{v})$ . The joint posterior and the full conditional densities cannot be expressed in closed form. Therefore, we cannot use the elementary Gibbs sampler and need to rely on a Metropolis-within-Gibbs strategy to approximate the posterior density. In this case, there is an R package called `bayesGARCH` that can be used to perform the analysis. Implementations details are given in ?. For this example, we use part of the CAC returns from the R dataset `EuStockMarkets`. The graphical output of the package is displayed in Figure 3.2, which shows the traces for each parameter (after burnin) and the corresponding posterior densities. We also use the R package `IDPmisc` to graph the pairwise results; see Figure 3.3.

### 3.5 Some proofs

Let  $(X, d)$  be a Polish space equipped with its Borel sigma-field  $\mathcal{X}$ . Let  $\theta : X^{\mathbb{N}} \rightarrow X^{\mathbb{N}}$  and  $\tilde{\theta} : X^{\mathbb{Z}} \rightarrow X^{\mathbb{Z}}$  be the shift operators defined by: for all  $\mathbf{x} = (x_t)_{t \in \mathbb{N}} \in X^{\mathbb{N}}$  and all  $\tilde{\mathbf{x}} = (\tilde{x}_t)_{t \in \mathbb{Z}} \in X^{\mathbb{Z}}$ ,

$$\theta(\mathbf{x}) = (y_t)_{t \in \mathbb{N}}, \quad \text{where } y_t = x_{t+1}, \quad \forall t \in \mathbb{N}, \quad (3.64)$$

$$\tilde{\theta}(\tilde{\mathbf{x}}) = (\tilde{y}_t)_{t \in \mathbb{Z}}, \quad \text{where } \tilde{y}_t = \tilde{x}_{t+1}, \quad \forall t \in \mathbb{Z}. \quad (3.65)$$

Assume that  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}}, \mathbb{P}, \tilde{\theta})$  is a measure-preserving ergodic dynamical system. Denote by  $\mathbb{E}$  the expectation operator associated to  $\mathbb{P}$ .

Let  $(\bar{\ell}^\theta, \theta \in \Theta)$  be a family of measurable functions,  $\bar{\ell}^\theta : X^{\mathbb{Z}} \rightarrow \mathbb{R}$ , indexed by  $\theta \in \Theta$  where  $(\Theta, d)$  is a compact metric space and denote  $\bar{L}_n^\theta := n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k$ . Moreover, consider  $(L_n^\theta, n \in \mathbb{N}^*, \theta \in \Theta)$  a family of upper-semicontinuous functions  $L_n^\theta : X^{\mathbb{Z}} \rightarrow \mathbb{R}$  indexed by  $n \in \mathbb{N}^*$  and  $\theta \in \Theta$ . Consider the following assumptions:

$$\mathbf{H} \text{ 3.39 } \mathbb{E}[\sup_{\theta \in \Theta} \bar{\ell}_+^\theta] < \infty,$$

$$\mathbf{H} \text{ 3.40 } \mathbb{P} - \text{a.s.}, \text{ the function } \theta \mapsto \bar{\ell}^\theta \text{ is upper-semicontinuous,}$$

$$\mathbf{H} \text{ 3.41 } \lim_{n \rightarrow \infty} \sup_{\theta \in \Theta} |L_n^\theta - \bar{L}_n^\theta| = 0, \quad \mathbb{P} - \text{a.s.}$$

Let  $\{\bar{\theta}_n : n \in \mathbb{N}^*\} \subset \Theta$  and  $\{\hat{\theta}_n : n \in \mathbb{N}^*\} \subset \Theta$  such that for all  $n \geq 1$ ,

$$\bar{\theta}_n \in \arg \max_{\theta \in \Theta} \bar{L}_n^\theta, \quad \hat{\theta}_n \in \arg \max_{\theta \in \Theta} L_n^\theta.$$

The proofs are adapted from ?.

**Theorem 3.42.** Assume **H 3.39–H 3.40**. Then

$$\lim_{n \rightarrow \infty} d(\bar{\theta}_n, \Theta_\star) = 0, \quad \mathbb{P} - \text{a.s. where } \Theta_\star := \arg \max_{\theta \in \Theta} \mathbb{E}[\bar{\ell}^\theta]. \quad (3.66)$$

If, in addition, **H 3.41** holds, then  $\lim_{n \rightarrow \infty} d(\hat{\theta}_n, \Theta_\star) = 0, \quad \mathbb{P} - \text{a.s. and}$

$$\lim_{n \rightarrow \infty} L_n^{\hat{\theta}_n} = \sup_{\theta \in \Theta} \mathbb{E}[\bar{\ell}^\theta], \quad \mathbb{P} - \text{a.s.} \quad (3.67)$$

$$\text{for all } \theta \in \Theta, \quad \lim_{n \rightarrow \infty} L_n^\theta = \mathbb{E}[\bar{\ell}^\theta], \quad \mathbb{P} - \text{a.s.} \quad (3.68)$$

PROOF. First note that according to the Birkhoff ergodic theorem (??) and **H 3.39**, for all  $\theta \in \Theta$ ,  $\lim_{n \rightarrow \infty} \bar{L}_n^\theta$  exists  $\mathbb{P} - \text{a.s.}$ , and

$$\lim_{n \rightarrow \infty} \bar{L}_n^\theta = \lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k = \mathbb{E}[\bar{\ell}^\theta], \quad \mathbb{P} - \text{a.s.} \quad (3.69)$$

Let  $K$  be a compact subset of  $\Theta$ . For all  $\theta_0 \in K, \mathbb{P} - \text{a.s.}$ ,

$$\begin{aligned} & \limsup_{\rho \rightarrow 0} \limsup_{n \rightarrow \infty} \sup_{\theta \in B(\theta_0, \rho)} n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k \\ & \leq \limsup_{\rho \rightarrow 0} \limsup_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} \sup_{\theta \in B(\theta_0, \rho)} \bar{\ell}^\theta \circ \tilde{\theta}^k = \limsup_{\rho \rightarrow 0} \mathbb{E}[\sup_{\theta \in B(\theta_0, \rho)} \bar{\ell}^\theta], \end{aligned} \quad (3.70)$$

where the last equality follows from **H** 3.39 and the Birkhoff ergodic theorem (??). Moreover, by the monotone convergence theorem applied to the nondecreasing function  $\rho \mapsto \sup_{\theta \in B(\theta_0, \rho)} \bar{\ell}^\theta$ , we have

$$\limsup_{\rho \rightarrow 0} \mathbb{E}[\sup_{\theta \in B(\theta_0, \rho)} \bar{\ell}^\theta] = \mathbb{E}[\limsup_{\rho \rightarrow 0} \sup_{\theta \in B(\theta_0, \rho)} \bar{\ell}^\theta] \leq \mathbb{E}[\bar{\ell}^{\theta_0}], \quad (3.71)$$

where the last inequality follows from **H** 3.40. Combining (3.70) and (3.71), we obtain that for all  $\eta > 0$  and  $\theta_0 \in K$ , there exists  $\rho^{\theta_0} > 0$  satisfying

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in B(\theta_0, \rho^{\theta_0})} n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k \leq \mathbb{E}[\bar{\ell}^{\theta_0}] + \eta \leq \sup_{\theta \in K} \mathbb{E}[\bar{\ell}^\theta] + \eta, \quad \mathbb{P} - \text{a.s.}$$

Since  $K$  is a compact subset of  $\Theta$ , we can extract a finite subcover of  $K$  from  $\bigcup_{\theta_0 \in K} B(\theta_0, \rho^{\theta_0})$ , so that

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in K} n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k \leq \sup_{\theta \in K} \mathbb{E}[\bar{\ell}^\theta] + \eta, \quad \mathbb{P} - \text{a.s.} \quad (3.72)$$

Since  $\eta$  is arbitrary, we obtain

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in K} n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k \leq \sup_{\theta \in K} \mathbb{E}[\bar{\ell}^\theta], \quad \mathbb{P} - \text{a.s.} \quad (3.73)$$

Moreover, (3.71) implies

$$\limsup_{\rho \rightarrow 0} \sup_{\theta \in B(\theta_0, \rho)} \mathbb{E}[\bar{\ell}^\theta] \leq \limsup_{\rho \rightarrow 0} \mathbb{E}[\sup_{\theta \in B(\theta_0, \rho)} \bar{\ell}^\theta] \leq \mathbb{E}[\bar{\ell}^{\theta_0}],$$

This shows that  $\theta \mapsto \mathbb{E}[\bar{\ell}^\theta]$  is upper-semicontinuous. As a consequence,  $\Theta_* := \arg \max_{\theta \in \Theta} \mathbb{E}[\bar{\ell}^\theta]$  is a closed and nonempty subset of  $\Theta$  and therefore, for all  $\varepsilon > 0$ ,  $K_\varepsilon := \{\theta \in \Theta; d(\theta, \Theta_*) \geq \varepsilon\}$  is a compact subset of  $\Theta$ . Using again the upper-semicontinuity of  $\theta \mapsto \mathbb{E}[\bar{\ell}^\theta]$ , there exists  $\theta_\varepsilon \in K_\varepsilon$  such that for all  $\theta_* \in \Theta_*$ ,

$$\sup_{\theta \in K_\varepsilon} \mathbb{E}[\bar{\ell}^\theta] = \mathbb{E}[\bar{\ell}^{\theta_\varepsilon}] < \mathbb{E}[\bar{\ell}^{\theta_*}].$$

Finally, combining this inequality with (3.73), we obtain that  $\mathbb{P} - \text{a.s.}$ ,

$$\begin{aligned} \limsup_{n \rightarrow \infty} \sup_{\theta \in K_\varepsilon} \bar{\Gamma}_n^\theta &= \limsup_{n \rightarrow \infty} \sup_{\theta \in K_\varepsilon} n^{-1} \sum_{k=0}^{n-1} \bar{\ell}^\theta \circ \tilde{\theta}^k \leq \sup_{\theta \in K_\varepsilon} \mathbb{E}[\bar{\ell}^\theta] \\ &< \mathbb{E}[\bar{\ell}^{\theta_*}] \stackrel{(1)}{=} \lim_{n \rightarrow \infty} \bar{\Gamma}_n^{\theta_*} \leq \liminf_{n \rightarrow \infty} \bar{\Gamma}_n^{\tilde{\theta}_n}, \end{aligned} \quad (3.74)$$

where (1) follows from (3.69). This inequality ensures that  $\tilde{\theta}_n \notin K_\varepsilon$  for all  $n$  larger to some  $\mathbb{P} - \text{a.s.}$  finite integer-valued random variable. This completes the proof of (3.66) since  $\varepsilon$  is arbitrary.

Eq. (3.68) follows from (3.69) and **H** 3.41. Let  $\theta_*$  be any point in  $\Theta_*$ . Then,  $\mathbb{P} - \text{a.s.}$ ,

$$\begin{aligned} \mathbb{E}[\bar{\ell}^{\theta_*}] &\stackrel{(1)}{=} \liminf_{n \rightarrow \infty} \bar{\Gamma}_n^{\theta_*} \stackrel{(2)}{\leq} \liminf_{n \rightarrow \infty} \bar{\Gamma}_n^{\tilde{\theta}_n} \leq \limsup_{n \rightarrow \infty} \bar{\Gamma}_n^{\tilde{\theta}_n} \\ &= \limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} \bar{\Gamma}_n^\theta \stackrel{(3)}{\leq} \sup_{\theta \in \Theta} \mathbb{E}[\bar{\ell}^\theta] = \mathbb{E}[\bar{\ell}^{\theta_*}], \end{aligned}$$

where (1) follows from (3.69), (2) follows from the definition of  $\tilde{\theta}_n$  and (3) is obtained by applying (3.73) with  $K = \Theta$ . Thus,

$$\lim_{n \rightarrow \infty} \bar{\Gamma}_n^{\tilde{\theta}_n} = \mathbb{E}[\bar{\ell}^{\theta_*}], \quad \mathbb{P} - \text{a.s.} \quad (3.75)$$

Denote  $\delta_n := \sup_{\theta \in \Theta} |\bar{\Gamma}_n^\theta - \bar{\Gamma}_n^{\theta_*}|$ . We get

$$\bar{\Gamma}_n^{\tilde{\theta}_n} - \delta_n \stackrel{(1)}{\leq} \bar{\Gamma}_n^{\tilde{\theta}_n} \stackrel{(2)}{\leq} \bar{\Gamma}_n^{\hat{\theta}_n} \stackrel{(1)}{\leq} \bar{\Gamma}_n^{\hat{\theta}_n} + \delta_n \stackrel{(3)}{\leq} \bar{\Gamma}_n^{\hat{\theta}_n} + \delta_n. \quad (3.76)$$

where (1) follows from the definition of  $\delta_n$ , (2) from the definition of  $\hat{\theta}_n$  and (3) from the definition of  $\tilde{\theta}_n$ . Combining the above inequalities with (3.75) and **H** 3.41 yields (3.67). (3.76) also implies that

$$\lim_{n \rightarrow \infty} \bar{\Gamma}_n^{\hat{\theta}_n} = \mathbb{E}[\bar{\ell}^{\theta_*}], \quad \mathbb{P} - \text{a.s.}$$

which yields, using (3.74),

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in K_\varepsilon} \bar{L}_n^\theta < \liminf_{n \rightarrow \infty} \bar{L}_n^{\hat{\theta}_n} = \limsup_{n \rightarrow \infty} \bar{L}_n^{\hat{\theta}_n} = \mathbb{E}[\bar{\ell}^{\theta_*}], \quad \mathbb{P} - \text{a.s.}$$

where  $K_\varepsilon := \{\theta \in \Theta; d(\theta, \Theta_*) \geq \varepsilon\}$ . Therefore,  $\hat{\theta}_n \notin K_\varepsilon$  for all  $n$  larger to some  $\mathbb{P} - \text{a.s.}$ -finite integer-valued random variable. The proof is completed since  $\varepsilon$  is arbitrary.  $\blacksquare$

## Exercises

**Exercise 3.43.** Let  $\{X_t, t \in \mathbb{N}\}$  be an  $X$ -valued stochastic process, where  $(X, \mathcal{X})$  is a measurable space. Assume that for each positive integer  $n$ , the random vector  $(X_0, \dots, X_{n-1})$  has a joint density  $x_{0:n-1} \mapsto p^\theta(x_{0:n-1})$  with respect to the product measure  $\lambda^{\otimes n}$  on the product space  $(X^n, \mathcal{X}^{\otimes(n-1)})$ , where  $\lambda \in \mathbb{M}_1(\mathcal{X})$ . For simplicity, it is assumed that for all  $n \in \mathbb{N}$  and all  $x_{0:n-1}$ ,  $p^\theta(x_{0:n-1}) > 0$ . Denote by  $p^\theta(x_n | x_{0:n-1})$  the conditional density of  $X_n$  given  $X_{0:n-1}$ ,  $p^\theta(x_n | x_{0:n-1}) = p^\theta(x_{0:n}) / p^\theta(x_{0:n-1})$ . Assume that, for all  $n \in \mathbb{N}$ , and  $x_{0:n-1} \in X^n$ ,  $\theta \mapsto p^\theta(x_{0:n-1})$  is differentiable and that

$$\nabla p^\theta(x_{0:n-1}) = \int \nabla p^\theta(x_{0:n}) \lambda(dx_n).$$

Prove that the score function (the gradient of the log-likelihood)  $\sum_{k=1}^n \nabla \log p^\theta(X_k | X_{0:k-1})$  is a martingale adapted to the natural filtration of the process.

**Exercise 3.44 (Conditionally Gaussian MLE for an AR(1)).** Consider a causal AR(1) model,  $X_t = \phi X_{t-1} + \sigma Z_t$ , where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise having a density  $f$ , which is symmetric around the origin and satisfies  $\int_{-\infty}^{\infty} x^2 f(x) \text{Leb}(dx) < \infty$ . We put  $\theta = (\phi, \sigma) \in \Theta$ , where  $\Theta$  is a compact subset of  $(-1, 1) \times \mathbb{R}_*^+$ . We denote by  $\theta_* = (\phi_*, \sigma_*^2) \in \Theta^o$  the true value of the parameters. We observe  $(X_0, \dots, X_n)$  and we are interested in estimating these parameters. We fit the model using the conditional *Gaussian* maximum likelihood estimator. We denote by  $(\hat{\phi}_n, \hat{\sigma}_n^2)$  these estimators.

- Derive the expressions of  $(\hat{\phi}_n, \hat{\sigma}_n^2)$ .
- Show that this estimator is strongly consistent.
- Show that the estimator  $\hat{\phi}_n$  is asymptotically normal. Compute the asymptotic variance of this estimator.
- Assume that  $f$  is  $t$ -distribution with  $\nu > 2$  degrees of freedom. Compute the variance of  $\hat{\phi}_n$  (as a function of  $\nu$ ) and construct confidence intervals. Discuss the result.
- Assume that  $\int_{-\infty}^{\infty} x^4 f(x) \text{Leb}(dx) < \infty$ . Show that  $(\hat{\phi}_n, \hat{\sigma}_n^2)$  is asymptotically normal.
- Determine confidence regions for  $(\hat{\phi}_n, \hat{\sigma}_n^2)$  when  $f$  is a Student  $t$ -distribution with  $\nu > 4$  degrees of freedom.

**Exercise 3.45 (Conditionally Gaussian MLE for an AR(1)).** We use the same notations and definitions as in Exercise 3.44. We now fit the model using the conditional maximum likelihood estimator. We assume in this exercise that the model is well specified. We denote by  $(\tilde{\phi}_n, \tilde{\sigma}_n^2)$  these estimators.

- Derive the estimating equations for  $(\tilde{\phi}_n, \tilde{\sigma}_n^2)$ .
- Derive explicitly the likelihood equations when  $f$  is a  $t$ -distribution with  $\nu$ -degrees of freedom.
- Show that these estimators (assuming that the model is well-specified) are consistent and asymptotically normal.
- Compute the asymptotic variance of the resulting estimator.
- Compare the asymptotic variance with Exercise 3.44.

**Exercise 3.46.** Verify (3.32).

**Exercise 3.47 (AR process with nonnegative innovation).** Consider the nonnegative autoregressive process  $X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + Z_t$ , where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise with left endpoint of their common distributions being zero and  $(\phi_1, \dots, \phi_p) \in (0, 1)^p$ . This type of model is suitable for modelling

applications of data that are inherently nonnegative, such as streams flow in hydrology, interarrival times, etc. Based on the observations,  $(X_0, \dots, X_n)$ , we are interested in estimating the parameters. The parameters  $(\phi_1, \dots, \phi_p)$  are estimated using the maximum likelihood estimator associated to the assumption that the common distribution of the  $Z$ 's is unit exponential, so that  $\mathbb{P}(Z_1 > x) = e^{-x}$ ,  $x \geq 0$ . See ?.

- (a) Consider first the case  $p = 1$ . Show that the conditional likelihood function is given (up to a constant) by

$$\prod_{t=1}^n \mathbb{1}_{\{X_t - \phi_1 X_{t-1} \geq 0\}} \exp \left( \phi_1 \sum_{t=1}^n X_{t-1} \right).$$

- (b) Show that the conditional maximum likelihood estimator is given by

$$\hat{\phi}_{1,n} = \min_{1 \leq t \leq n} (X_t / X_{t-1}).$$

- (c) Consider now the general autoregressive case. Show that the likelihood function is given (up to a constant) by

$$\prod_{t=p}^n \mathbb{1}_{\{X_t - \phi_1 X_{t-1} \geq 0\}} \exp \left( \phi_1 \sum_{t=p}^n X_{t-1} + \dots + \phi_p \sum_{t=p}^n X_{t-p} \right).$$

- (d) Argue that the maximum likelihood estimator is approximately determined by solving the linear program

$$\max \left( \sum_{i=1}^p \phi_i \right) \quad \text{subject to} \quad X_t \geq \sum_{i=1}^p \phi_i X_{t-i} \text{ for } t \in \{1, \dots, n\}.$$

Determining the asymptotic properties of these estimators is difficult (see ?); the rate of convergence can be in such a case faster than  $n^{1/2}$ .

**Exercise 3.48 (Proof of Theorem 3.18).**

- (a) Show that

$$\begin{aligned} n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\hat{\theta}_n}(X_{t-p:t-1}; X_t) &= 0 = n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\theta_*}(X_{t-p:t-1}; X_t) \\ &\quad + n^{-1} \sum_{t=p}^n \left( \int_0^1 \nabla^2 \ln q^{\theta_{n,s}}(X_{t-p:t-1}; X_t) ds \right) \sqrt{n}(\hat{\theta}_n - \theta_*), \end{aligned} \quad (3.77)$$

where  $\theta_{n,s} = s\hat{\theta}_n + (1-s)\theta_*$ .

- (b) Show that

$$n^{-1/2} \sum_{t=p}^n \nabla \ln q^{\theta_*}(X_{t-p:t-1}; X_t) \xrightarrow{\mathcal{L}_{\mathbb{P}^{\theta_*}}} \mathbf{N}(0, \mathcal{J}(\theta_*)).$$

- (c) Show that, for all  $i, j \in \{1, \dots, d\}$ ,

$$n^{-1} \sum_{t=p}^n \int_0^1 \frac{\partial^2 \ln q^{\theta_{n,s}}}{\partial \theta_i \partial \theta_j}(X_{t-p:t-1}; X_t) ds \xrightarrow{\mathbb{P}^{\theta_*-\text{prob}}} \mathbb{E}^{\theta_*} \left[ \frac{\partial^2 \ln q^{\theta_*}}{\partial \theta_i \partial \theta_j}(X_{0:p-1}; X_p) \right].$$

- (d) Show that

$$n^{-1} \sum_{t=p}^n \left( \int_0^1 \nabla^2 \ln q^{\theta_{n,s}}(X_{t-p:t-1}; X_t) ds \right) \xrightarrow{\mathbb{P}^{\theta_*-\text{prob}}} \mathcal{J}(\theta_*).$$

- (e) Conclude.

**Exercise 3.49 (Autoregressive model with ARCH error).** Consider the  $p$ th order autoregressive model with ARCH(1) error

$$Y_t = \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \dots + \alpha_p Y_{t-p} + \varepsilon_t, \quad (3.78)$$

where  $\varepsilon_t | \mathcal{F}_{t-1} \sim N(0, \beta_0 + \beta_1 \varepsilon_{t-1}^2)$  and  $\{\mathcal{F}_t, t \in \mathbb{N}\}$  is the natural filtration of the process  $\{Y_t, t \in \mathbb{N}\}$ , i.e., for any nonnegative integer  $t$ ,  $\mathcal{F}_t = \sigma(Y_s, s \leq t)$ . It is assumed that  $\beta_0$  is strictly positive and  $0 \leq \beta_1 < 1$ .

- (a) Show that the unique strict-sense stationary solution to the recursion  $\varepsilon_t = \sqrt{\beta_0 + \beta_1 \varepsilon_{t-1}^2} Z_t$ , where  $\{Z_t, t \in \mathbb{Z}\}$  is a strong Gaussian noise with zero-mean and unit-variance is given by

$$\varepsilon_t = \beta_0 \sum_{l=0}^{\infty} \beta_1^l \left( \prod_{i=0}^l Z_{t-i}^2 \right). \quad (3.79)$$

- (b) Let  $r$  be a nonnegative integer. Show that  $m_{2r} = \mathbb{E}[\varepsilon_t^{2r}]$  exists if and only if  $\eta_r = \beta_1^r \prod_{j=1}^r (2j-1) < 1$ .  
(c) Show that if  $3\beta_1^2 < 1$ , then  $\text{Var}(\varepsilon_t^2) = 2\sigma^4(1 - 3\beta_1^2)^{-1}$ , and  $\text{Cov}(\varepsilon_t^2, \varepsilon_{t-j}^2) = \beta_1^j \text{Var}(\varepsilon_t^2)$  where  $\sigma^2 = \beta_0(1 - \beta_1)^{-1}$  is the stationary unconditional variance.  
(d) Show that, if  $\eta_r < 1$ , then

$$n^{-1} \sum_{t=1}^n \varepsilon_t^{2r} \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}[\varepsilon_t^{2r}], \quad \text{and} \quad n^{-1} \sum_{t=1}^n \varepsilon_t^{2r-1} \xrightarrow{\mathbb{P}\text{-a.s.}} 0.$$

- (e) Show that if  $3\beta_1^2 < 1$ , then for a fixed  $j \neq 0$ ,

$$n^{-1} \sum_{t=1}^n \varepsilon_t \varepsilon_{t-j} \xrightarrow{\mathbb{P}\text{-a.s.}} 0 \quad \text{and} \quad n^{-1} \sum_{t=1}^n \varepsilon_t^2 \varepsilon_{t-j}^2 \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}[\varepsilon_t^2 \varepsilon_{t-j}^2].$$

- (f) Given  $(Y_0, \dots, Y_n)$  we want to estimate  $\alpha = (\alpha_1, \dots, \alpha_p)$  and  $\beta = (\beta_0, \beta_1)$ . We write  $\theta = (\alpha, \beta)$ . Show that the log-likelihood conditional on  $(Y_p, Y_{p-1}, \dots, Y_0)$  is given by

$$\bar{\mathcal{L}}_n^\theta(Y_{p:n} | Y_{0:p-1}) = (n-p)^{-1} \sum_{t=p+1}^n \ln p^\theta(Y_t | Y_{t-1:t-p-1})$$

with

$$p^\theta(Y_t | Y_{t-1:t-p-1}) = (2\pi h_t(\theta))^{-1/2} \exp[-(Y_t - \mathbf{X}'_{t-1} \alpha)^2 / 2h_t(\theta)],$$

$$h_t(\theta) = \beta_0 + \beta_1 (Y_{t-1} - \mathbf{X}'_{t-2} \alpha)^2,$$

where  $\mathbf{X}_t = (Y_t, Y_{t-1}, \dots, Y_{t-p+1})$ .

- (g) We assume in the sequel that the model is well-specified. State sufficient conditions upon which the MLE estimator is strongly consistent.  
(h) State sufficient conditions upon which the MLE estimator is asymptotically normal, i.e.,  $n^{1/2}(\hat{\theta}_n - \theta_*) \xrightarrow{\mathcal{L}_\mathbb{P}} N(0, H^{-1}(\theta_*))$ .  
(i) Show that the asymptotic covariance matrix  $H(\theta_*) = (H_{i,j}(\theta_*))_{1 \leq i,j \leq 2}$  is the  $(p+2) \times (p+2)$  block matrix with  $H_{1,2}(\theta_*) = H_{2,1}(\theta_*) = 0$ ,

$$H_{1,1}(\theta_*) = \mathbb{E}^{\theta_*}[v_t^{-1}(\theta_*) \mathbf{X}'_{t-1} \mathbf{X}_{t-1}] + 2\beta_{*,1} \mathbb{E}^{\theta_*}[v_t^{-1}(\theta_*) \mathbf{X}'_{t-2} \mathbf{X}_{t-2}]$$

$$- 2\beta_{*,1} \beta_{*,0} \mathbb{E}^{\theta_*}[v_t^{-2}(\theta_*) \mathbf{X}'_{t-2} \mathbf{X}_{t-2}],$$

where  $\{v_t(\theta_*), t \in \mathbb{Z}\}$  is the unique strict sense stationary non anticipative solution of  $v_t(\theta_*) = \beta_{*,0} + \beta_{*,1} v_{t-1}(\theta_*) Z_{t-1}^2$ . Compute similarly  $H_{2,2}(\theta_*)$ .

- (j) Show that the matrix  $H$  may be estimated consistently by replacing the parameters by their estimates and expectations by sample averages in the above expression.  
(k) Construct asymptotic confidence regions for the parameters.

**Exercise 3.50 (EXPAR(1) model).** Consider the EXPAR(1) model (see ??)

$$X_t = \{\phi_0 + \phi_1 \exp(-\gamma X_{t-1}^2)\} X_{t-1} + \sigma Z_t, \quad (3.80)$$



where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise with zero mean, unit variance independent of  $X_0$ . Assume that the distribution of  $Z_0$  has a continuous density  $f$  with respect to the Lebesgue measure, which is everywhere positive and  $\mathbb{E}[Z_0^{2q}] < \infty$  for some  $q \in \mathbb{N}_*$ . Denote by  $\theta = (\phi_0, \phi_1, \gamma, \sigma^2) \in \Theta$  the unknown parameters, where  $\Theta$  is a compact subset of

$$\{(\phi_0, \phi_1, \gamma, \sigma) : \gamma > 0, \sigma^2 > 0, \max(|\phi_0|, |\phi_0 + \phi_1|) < 1\}.$$

For  $\theta \in \Theta$ , denote  $a^\theta(x) = \phi_0 + \phi_1 \exp(-\gamma x^2)$  and  $V_q(x) = 1 + |x|^{2q}$ .

- (a) Show that the Markov chain is  $V_q$ -geometrically ergodic.
- (b) We observe  $(X_0, X_1, \dots, X_n)$  ( $n+1$ ) observations from the strict sense stationary solution of (3.80) associated to the “true” parameters  $\theta_* \in \Theta$ . We fit the parameters of the model using the Gaussian likelihood. Show that for all  $(x_0, x_1) \in \mathbb{R}^2$ ,

$$\ln q^\theta(x_0, x_1) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (x_1 - a^\theta(x_0))^2.$$

- (c) Assume that  $q = 1$ . Show that the sequence of estimators  $\hat{\theta}_n \rightarrow \mathbb{P}_{\theta_*}$ -a.s.  $\theta_*$ .
- (d) State conditions upon which the sequence of estimators  $\{\hat{\theta}_n, n \in \mathbb{N}\}$  is asymptotically normal and compute the asymptotic covariance matrix.

**Exercise 3.51 (Nonlinear autoregression).** Consider the NLAR(1) process

$$X_t = f^\theta(X_{t-1}) + Z_t, \quad (3.81)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is a strong white noise independent of  $X_0$  and  $\theta \in \Theta$  a compact subset of  $\mathbb{R}$ . Assume that for all  $\theta \in \Theta$ ,  $\lim_{|x| \rightarrow \infty} |f^\theta(x)|/|x| < 1$ ,  $f^\theta : \mathbb{R} \rightarrow \mathbb{R}$  is a twice continuously differentiable function, the distribution of  $Z_0$  has a continuous density with respect to the Lebesgue measure, which is everywhere positive and that  $\mathbb{E}[Z_0^4] < \infty$ .

- (a) Show that the conditional Gaussian likelihood is given by

$$\ln p^\theta(X_{p:n}|X_{0:p-1}) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{t=1}^n \{X_t - f^\theta(X_{t-1})\}^2.$$

- (b) Assume  $\{X_t, t \in \mathbb{N}\}$  is strictly stationary such that  $\mathbb{E}[\sup_{\theta \in \Theta} (X_1 - f^\theta(X_0))^2] < \infty$ . Then, show that the sequence of MLE estimators  $\hat{\theta}_n$  is strongly consistent.
- (c) Assume that the model is well-specified and denote by  $\theta_* \in \Theta$  the “true” values of the unknown parameters. Show that the model is identifiable if the condition  $f^\theta(X_0) = f^{\theta_*}(X_0)$   $\mathbb{P}^{\theta_*}$ -a.s. implies that  $\theta = \theta_*$ .
- (d) Give conditions upon which the sequence of estimators is asymptotically normal, and give the expressions of the asymptotic covariance matrix both in the well-specified and in the misspecified cases.

**Exercise 3.52 (Proof of Theorem 3.28).**

- (a) Show that, under **H** 3.27, for all  $\theta \in \Theta$ ,

$$\begin{aligned} & \mathbb{E}[\bar{\ell}^\theta(Y_{-\infty:1})] \\ &= \mathbb{E}\left[\mathbb{E}\left[\ln q(f_{Y_{-\infty:0}}^\theta, Y_1) \mid Y_s, s \leq 0\right]\right] = \mathbb{E}\left[\int \mathcal{Q}^*(f_{Y_{-\infty:0}}^{\theta_*}, dy) \ln q(f_{Y_{-\infty:0}}^\theta, y)\right] \\ &\leq \mathbb{E}\left[\int \mathcal{Q}^*(f_{Y_{-\infty:0}}^{\theta_*}, dy) \ln q(f_{Y_{-\infty:0}}^{\theta_*}, y)\right] = \mathbb{E}[\bar{\ell}^{\theta_*}(Y_{-\infty:1})]. \end{aligned} \quad (3.82)$$

- (b) Show that (3.45) implies that  $\theta = \theta_*$ .
- (c) Conclude.

**Exercise 3.53 (Proof of Lemma 3.34).** Let  $\mathcal{F}$  be the filtration  $\mathcal{F} = (\mathcal{F}_n)_{n \in \mathbb{N}}$  where  $\mathcal{F}_n = \sigma(Y_s, s \leq n)$  and let

$$M_n := \sum_{t=1}^n \nabla \ln q(f_{Y_{-\infty:t-1}}^{\theta_*}, Y_t) = \sum_{t=1}^n \varphi(\nabla f_{Y_{-\infty:t-1}}^{\theta_*}, f_{Y_{1:t-1}}^{\theta_*}(x), Y_t) .$$

(a) Show that  $\mathbb{E}(\|M_n\|^2) < \infty$  and

$$\begin{aligned} \mathbb{E}^{\theta_*} \left[ \nabla \ln q(f_{Y_{-\infty:t-1}}^{\theta_*}, Y_t) \middle| \mathcal{F}_{t-1} \right] \\ = \nabla f_{Y_{-\infty:t-1}}^{\theta_*} \int \mathcal{Q}^*(f_{Y_{-\infty:t-1}}^{\theta_*}, dy) \frac{\partial \ln q(f_{Y_{-\infty:t-1}}^{\theta_*}, y)}{\partial x} = 0 . \end{aligned}$$

(b) Show that  $\{M_t, t \geq 1\}$  is a square integrable  $\mathcal{F}$ -martingale with stationary and ergodic increments.  
(c) Conclude.

**Exercise 3.54 (Proof of Lemma 3.35).** Denote by  $A_t(\theta) = \frac{\partial^2 \ln q}{\partial \theta_i \partial \theta_j}(f_{Y_{-\infty:t-1}}^{\theta}, Y_t)$ .

(a) Show that  $n^{-1} \sum_{t=1}^n A_t(\theta_*) \rightarrow_{\mathbb{P}} \mathbb{E}[A(\theta_*)]$ .  
(b) Let  $\varepsilon > 0$ . Show that we can choose  $0 < \eta < \rho$  such that

$$\mathbb{E} \left( \sup_{\theta \in B(\theta_*, \eta)} |A_1(\theta_*) - A_1(\theta)| \right) < \varepsilon . \quad (3.83)$$

(c) Show that

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{P} \left( n^{-1} \sum_{t=1}^n |A_t(\theta_*) - A_t(\theta_n)| \geq \varepsilon, \theta_n \in B(\theta_*, \eta) \right) \\ \leq \limsup_{n \rightarrow \infty} \mathbb{P} \left( n^{-1} \sum_{t=1}^n \sup_{\theta \in B(\theta_*, \eta)} |A_t(\theta_*) - A_t(\theta)| \geq \varepsilon \right) = 0 , \end{aligned}$$

(d) Show that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left( n^{-1} \sum_{t=1}^n |A_t(\theta_*) - A_t(\theta_n)| \geq \varepsilon \right) = 0 .$$

(e) Show that  $n^{-1} \sum_{t=1}^n |A_t(\theta_*) - A_t(\theta_n)| \xrightarrow{\mathbb{P}\text{-prob}} 0$  and conclude.

**Exercise 3.55 (Conditional least-squares estimators).** Let  $p$  be an integer and  $\{P^\theta : \theta = (\varphi, \gamma) \in \Theta\}$  be a parametric family of Markov kernels on  $X^p \times \mathcal{X}$ , where  $\Theta$  is a compact subset of  $\mathbb{R}^d$  (here,  $\varphi$  is the parameter of interest and  $\gamma$  is a nuisance). Assume that for any  $\theta \in \Theta$ ,  $P^\theta$  has a unique stationary distribution satisfying  $\mathbb{E}^\theta[|X_0|] < \infty$ . Let  $\theta_* \in \Theta^o$ . Assume that the vector of observations  $(X_0, X_1, \dots, X_n)$  is, for each  $n$  a realization of the stationary ergodic Markov chain  $\{X_t, t \in \mathbb{Z}\}$  with transition kernel  $P^{\theta_*}$ . Finally, assume that

- there exists  $a^\varphi(X_{p-1}, \dots, X_0)$  a version of the conditional expectation  $a^\varphi(X_{p-1}, \dots, X_0) = \mathbb{E}^\theta[X_p | X_0, \dots, X_{p-1}]$  with  $\theta = (\varphi, \gamma)$  satisfying

$$\mathbb{E}^{\theta_*} \left[ \sup_{\varphi \in \Theta_\varphi} |a^\varphi(X_{p-1}, \dots, X_0)|^2 \right] < \infty ,$$

where  $\Theta_\varphi = \{\varphi, (\varphi, \gamma) \in \Theta\}$  is the canonical projection of  $\Theta$ .

- for any  $(x_0, \dots, x_{p-1}) \in X^p$ , the function  $\varphi \mapsto a^\varphi(x_{p-1}, \dots, x_0)$  is continuous.

We estimate the parameter  $\varphi \in \Theta_\varphi$  by the *conditional least-squares*, which amounts to estimating the model by using the likelihood of

$$Y_t = a^\varphi(Y_{t-1}, \dots, Y_{t-p}) + \sigma^2 Z_t \quad (3.84)$$

where  $\{Z_t, t \in \mathbb{Z}\}$  is a strong white noise with zero-mean and unit-variance. Assume in addition that ?? is satisfied and  $\mathbb{E}[Z_0^{2q}] < \infty$  for some  $q \in \mathbb{N}_*$ .

(a) Show that the log-likelihood for the misspecified model (3.84) is given by

$$\ln q^\vartheta(x_{p-1:0}; X_p) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \{x_p - a^\vartheta(x_{p-1}, \dots, x_0)\}^2.$$

- (b) Show that if  $a^{\vartheta_*}(X_{p-1}, \dots, X_0) = a^\vartheta(X_{p-1}, \dots, X_0) \mathbb{P}_{\theta_*} - \text{a.s.}$  implies that  $\vartheta = \vartheta_*$ , then the conditional least square estimator  $\{\hat{\vartheta}(n), n \in \mathbb{N}\}$  is strongly consistent.
- (c) State conditions upon which the conditional least-squares estimator is asymptotically normal and compute its asymptotic covariance matrix.

**Exercise 3.56.** Consider an EXPAR( $p$ ) model

$$X_t = \{a_1 + b_1 \exp(-\lambda_1 X_{t-d}^2)\} X_{t-1} + \dots + \{a_p + b_p \exp(-\lambda_p X_{t-p}^2)\} X_{t-p} + \sigma Z_t, \quad (3.85)$$

where  $\{Z_t, t \in \mathbb{Z}\}$  is a strong white noise with zero-mean and unit-variance. Assume in addition that ?? is satisfied and  $\mathbb{E}[Z_0^{2q}] < \infty$  for some  $q \in \mathbb{N}_*$ . We denote by  $\theta = \{(a_i, b_i, \lambda_i), i = 1, \dots, p, d, \sigma^2\} \in \Theta$  the unknown parameters, where  $\Theta$  is a compact subset of the set  $\{(a_i, b_i, \lambda_i), i = 1, \dots, p, d, \sigma^2\}$  satisfying  $\sigma^2 > 0$ ,  $d \in \{0, \dots, p\}$ ,  $\lambda_i > 0$ ,  $i \in \{1, \dots, p\}$  and  $c(z) = 1 - c_1 z - \dots - c_p z^p$  where  $c_i = \max(|a_i|, |a_i + b_i|)$ . For  $\theta \in \Theta$ , the Markov chain  $\mathbf{X}_t = [X_t, \dots, X_{t-p+1}]'$  given by (3.85) is  $V_q$ -geometrically ergodic with  $V_q(\mathbf{x}) = 1 + \|\mathbf{x}\|^q$ ; see ?? and ?. Denote by  $\theta_* \in \Theta$  the true parameter vector. Given the observations  $(X_0, \dots, X_n)$ , we estimate the parameters using the conditional least square (Gaussian likelihood)

$$-\frac{n-p}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=p}^n (X_t - a^\theta(X_{t-1}, \dots, X_{t-p}))^2,$$

where

$$a^\theta(x_{p-1}, \dots, x_0) = \sum_{i=1}^p \{a_i + b_i \exp(-\lambda_i x_{p-i}^2)\} x_{p-i}.$$

- (a) Compute the conditional least squares estimator of  $\theta$ .
- (b) Assume that  $(X_0, X_1, \dots, X_n)$  is a sample record of  $\{X_t, t \in \mathbb{N}\}$  which is the ergodic solution of (3.85) associated to the true parameter vector  $\theta_*$ . Show that if  $q = 1$ , then  $\hat{\theta}_n$  is strongly consistent.
- (c) Show that if  $q = 2$ , this estimator is asymptotically normal and compute the limiting covariance.

**Exercise 3.57 (First order RCA; ?? cont.).** We use the notations of ?. We denote by  $\theta = (\phi, \lambda, \sigma)$  the unknown parameters. We assume that the noise  $Z_t$  cannot take on only two values asymptotically.

- (a) Show that  $\mathbb{E}^\theta[X_t | \mathcal{F}_{t-1}^X] = \phi X_{t-1}$  and  $\text{Var}([ \theta ] X_t | \mathcal{F}_{t-1}^X) = \lambda^2 X_{t-1}^2 + \sigma^2$ ,  $\mathbb{P}^\theta - \text{a.s.}$
- (b) Show that if  $\{B_t, t \in \mathbb{Z}\}$  and  $\{Z_t, t \in \mathbb{Z}\}$  are jointly Gaussian, then the conditional distribution of  $X_t$  given  $\mathcal{F}_{t-1}^X$  is Gaussian with mean  $\mathbb{E}^\theta[X_t | \mathcal{F}_{t-1}^X]$  and variance  $\text{Var}([ \theta ] X_t | \mathcal{F}_{t-1}^X)$ . Write the likelihood of this model (referred to as in the sequel as the Gaussian likelihood).
- (c) Assume that  $\theta_* \in \Theta$ , where  $\Theta$  is a compact subset of

$$\{(\phi, \lambda, \sigma) \in \mathbb{R} \times \mathbb{R}_+^2 : \phi^2 + \lambda^2 < 1, \sigma > 0\}.$$

Assume that the observations  $(X_0, \dots, X_n)$  is a sample record from the strict-sense stationary solution  $\{X_t, t \in \mathbb{N}\}$  of (??) associated to some values of the parameters  $\theta_* \in \Theta$ . Prove that the Gaussian MLE  $\{\hat{\phi}_n\}_{n \geq 1}$  is a strongly consistent sequence of estimator of  $\phi_*$ .

- (d) Under which conditions is the Gaussian MLE  $\{(\hat{\lambda}_n, \hat{\sigma}_n^2)\}_{n \geq 0}$  a strongly consistent sequence of estimators of  $(\lambda, \sigma^2)$ ?
- (e) Assume that  $\theta_* \in \Theta^o$ , where  $\Theta$  is a compact subset of

$$\{(\phi, \lambda, \sigma) \in \mathbb{R} \times \mathbb{R}_+^2 : \phi^4 + 6\lambda^2\phi^2 + m_4\lambda^4 < 1, \sigma > 0\}.$$

Prove that the Gaussian MLE  $\{\hat{\phi}_n\}_{n \geq 0}$  is an asymptotically normal sequence of estimators of  $\phi_*$ . Compute a confidence interval for this parameter.

- (f) Under which conditions is the sequence  $\{(\hat{\lambda}_n, \hat{\sigma}_n^2)\}_{n \geq 0}$  asymptotically normal?

**Exercise 3.58.** Consider the stochastic process

$$X_t = A_t X_{t-1} + B_t, \quad (3.86)$$

where  $\{(A_t, B_t), t \in \mathbb{N}\}$  is a Gaussian white noise, independent of  $X_0$ , with mean zero and covariance matrix

$$\begin{pmatrix} \alpha^2 & \rho\alpha\beta \\ \rho\alpha\beta & \beta^2 \end{pmatrix}.$$

We denote by  $\theta = (\alpha^2, \beta^2, \rho) \in \Theta$  the unknown parameters, where  $\Theta$  is a compact subset of  $(0, 1) \times \mathbb{R}_*^+ \times (-1, 1)$ .

- Show that for any  $\theta \in \Theta$ , (3.86) admits a unique strict sense stationary solution, satisfying  $\mathbb{E}^\theta[X_0^2] < \infty$ .
- Find a necessary and sufficient condition on  $\theta$  under which (3.86) admits a unique strict sense stationary solution satisfying  $\mathbb{E}^\theta[X_0^4] < \infty$ .
- Write the conditional maximum likelihood estimator for the model (3.86).
- Write an algorithm in R to solve numerically this equation.
- Assume that the observation  $(X_0, \dots, X_n)$  is a realization of (3.86) for some parameter  $\theta_* \in \Theta$ . Prove that  $\hat{\theta}_n$  is a strongly consistent sequence of estimators of  $\theta_*$ .
- State conditions upon which this estimator is asymptotically normal and compute the asymptotic covariance matrix.

**Exercise 3.59 (APARCH(1,1); ?? cont.).** The APARCH(1,1) model of ? can be defined as follows:

$$Y_t = \sigma_t(\theta) \varepsilon_t \quad (3.87)$$

$$\sigma_t^\delta(\theta) = \alpha_0 + \alpha_1(Y_{t-1} - \gamma Y_{t-1})^\delta + \beta_1 \sigma_{t-1}^\delta(\theta), \quad (3.88)$$

where  $\{\varepsilon_t, t \in \mathbb{N}\}$  is a strong white Gaussian noise with zero-mean and unit variance and  $\theta = (\alpha_0, \alpha_1, \beta_1, \gamma, \delta) \in \Theta$  a compact set of

$$\{(\alpha_0, \alpha_1, \beta_1, \gamma, \delta) : \alpha_0 > 0, \alpha_1 > 0, 0 < \gamma < 1, \delta > 0\}.$$

The parameter  $\delta$ , ( $\delta > 0$ ) parameterizes a Box-Cox transformation of the conditional standard deviation  $\sigma_t(\theta)$ , while the parameters  $\gamma$  reflect the leverage effect.

- Show that the APARCH model is an observation-driven model. Determine the kernel  $q^\theta$  and the function  $f_y^\theta$ .
- Show that the maximum likelihood estimator is consistent.
- Compute the confidence intervals for the parameters.

Chapter

4

# Markov chain Monte Carlo methods

## Contents

<b>4.1</b>	<b>Metropolis-Hastings algorithm</b> .....	<b>66</b>
4.1.1	Metropolis-Hastings algorithms .....	66
4.1.2	Data augmentation .....	68
4.1.3	Two-stage Gibbs sampler .....	71
4.1.4	Hit-and-run algorithm .....	73
4.1.5	Gibbs sampling .....	74
<b>4.2</b>	<b>Ordering the asymptotic variances</b> .....	<b>76</b>

Let  $\pi$  be a given probability measure on  $(X, \mathcal{X})$ . The aim consists in constructing a Markov chain with invariant distribution  $\pi$ . To do so, we will construct Markov kernel  $P$  that are  $\pi$ -reversible. We first define  $\pi$ -reversibility and then show that if  $P$  is  $\pi$ -reversible then it admits  $\pi$  as invariant distribution.

**Definition 4.1.** Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  and  $\pi$  a probability measure on  $(X, \mathcal{X})$ . Then  $P$  is  $\pi$ -reversible iff for all  $h \in \mathbb{F}_+(X \times X)$ ,

$$\iint h(x, y) \pi(dx) Q(x, dy) = \iint h(x, y) \pi(dy) Q(y, dx)$$

Using infinitesimal notation (which is an abuse of notation), we can also say:  $P$  is  $\pi$ -reversible iff  $\pi(dx) Q(x, dy) = \pi(dy) Q(y, dx)$ .

**Lemma 4.2** *If  $P$  is reversible with respect to  $\pi$ , then  $\pi$  is a stationary distribution for  $P$ .*

PROOF. Setting  $f \equiv 1$ , we indeed obtain for any function  $g \in \mathbb{F}_+(X)$ ,

$$\iint \pi(dx) P(x, dy) g(y) = \iint \pi(dx) P(x, dy) g(x) = \int \pi(dx) g(x) .$$

■

## 4.1 Metropolis-Hastings algorithm

Markov chain Monte Carlo is a general method for the simulation of distributions known up to a multiplicative constant. Let  $\nu$  be a  $\sigma$ -finite measure on a state space  $(X, \mathcal{X})$  and let  $h_\pi \in \mathbb{F}_+(X)$  such that  $0 < \int_X h_\pi(x) \nu(dx) < \infty$ . Typically  $X$  is an open subset of  $\mathbb{R}^d$  and  $\nu$  is the Lebesgue measure or  $X$  is countable and  $\nu$  is the counting measure. This function is associated to a probability measure  $\pi$  on  $X$  defined by

$$\pi(A) = \frac{\int_A h_\pi(x) \nu(dx)}{\int_X h_\pi(x) \nu(dx)}. \quad (4.1)$$

We want to approximate expectations of functions  $f \in \mathbb{F}_+(X)$  with respect to  $\pi$

$$\pi(f) = \frac{\int_X f(x) h_\pi(x) \nu(dx)}{\int_X h_\pi(x) \nu(dx)}.$$

If the state space  $X$  is high-dimensional and  $h_\pi$  is complex, direct numerical integration is not an option. The classical Monte Carlo solution to this problem is to simulate i.i.d. random variables  $Z_0, Z_1, \dots, Z_{n-1}$  with distribution  $\pi$  and then to estimate  $\pi(f)$  by the sample mean

$$\hat{\pi}(f) = n^{-1} \sum_{i=0}^{n-1} f(Z_i). \quad (4.2)$$

This gives an unbiased estimate with standard deviation of order  $O(n^{-1/2})$  provided that  $\pi(f^2) < \infty$ . Furthermore, by the Central Limit Theorem, the normalized error  $\sqrt{n}(\hat{\pi}(f) - \pi(f))$  has a limiting normal distribution, so that confidence intervals are easily obtained.

The problem often encountered in applications is that it might be very difficult to simulate i.i.d. random variables with distribution  $\pi$ . Instead, the Markov chain Monte Carlo (MCMC) solution is to construct a Markov chain on  $X$  which has  $\pi$  as invariant probability. The hope is that regardless of the initial distribution  $\xi$ , the law of large numbers will hold, i.e.  $\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} f(X_k) = \pi(f)$   $\mathbb{P}_\xi$  - a.s.

At first sight, it may seem even more difficult to find such a Markov chain than to estimate  $\pi(f)$  directly. In the following subsections, we will exhibit several such constructions.

### 4.1.1 Metropolis-Hastings algorithms

Let  $Q$  be a Markov kernel having a density  $q$  with respect to  $\nu$  i.e.  $Q(x, A) = \int_A q(x, y) \nu(dy)$  for every  $x \in X$  and  $A \in \mathcal{X}$ .

The Metropolis-Hastings algorithm proceeds in the following way. An initial starting value  $X_0$  is chosen. Given  $X_k$ , a candidate move  $Y_{k+1}$  is sampled from  $Q(X_k, \cdot)$ . With probability  $\alpha(X_k, Y_{k+1})$ , it is accepted and the chain moves to  $X_{k+1} = Y_{k+1}$ . Otherwise the move is rejected and the chain remains at  $X_{k+1} = X_k$ . The probability  $\alpha(X_k, Y_{k+1})$  of accepting the move is given by

$$\alpha(x, y) = \begin{cases} \min\left(\frac{h_\pi(y)}{h_\pi(x)} \frac{q(y, x)}{q(x, y)}, 1\right) & \text{if } h_\pi(x)q(x, y) > 0, \\ 1 & \text{if } h_\pi(x)q(x, y) = 0. \end{cases} \quad (4.3)$$

The acceptance probability  $\alpha(x, y)$  only depends on the ratio  $h_\pi(y)/h_\pi(x)$ ; therefore, we only need to know  $h_\pi$  up to a normalizing constant. In Bayesian inference, this property plays a crucial role.

This procedure produces a Markov chain,  $\{X_k, k \in \mathbb{N}\}$ , with Markov kernel  $P$  given by

$$P(x, A) = \int_A \alpha(x, y) q(x, y) \nu(dy) + \bar{\alpha}(x) \delta_x(A), \quad (4.4)$$

with

$$\bar{\alpha}(x) = \int_{\mathbf{X}} \{1 - \alpha(x, y)\} q(x, y) \nu(dy) . \quad (4.5)$$

The quantity  $\bar{\alpha}(x)$  is the probability of remaining at the same point.

**Proposition 4.3** *The distribution  $\pi$  is reversible with respect to the Metropolis-Hastings kernel  $P$ .*

PROOF. Note first that for every  $x, y \in \mathbf{X}$ , it holds that

$$\begin{aligned} h_{\pi}(x) \alpha(x, y) q(x, y) &= \{h_{\pi}(x) q(x, y)\} \wedge \{h_{\pi}(y) q(y, x)\} \\ &= h_{\pi}(y) \alpha(y, x) q(y, x) . \end{aligned} \quad (4.6)$$

Thus for  $C \in \mathcal{X} \times \mathcal{X}$ ,

$$\begin{aligned} \iint h_{\pi}(x) \alpha(x, y) q(x, y) \mathbb{1}_C(x, y) \nu(dx) \nu(dy) \\ = \iint h_{\pi}(y) \alpha(y, x) q(y, x) \mathbb{1}_C(x, y) \nu(dx) \nu(dy) . \end{aligned} \quad (4.7)$$

On the other hand,

$$\begin{aligned} \iint h_{\pi}(x) \delta_x(dy) \bar{\alpha}(x) \mathbb{1}_C(x, y) \nu(dx) \\ = \int h_{\pi}(x) \bar{\alpha}(x) \mathbb{1}_C(x, x) \nu(dx) = \int h_{\pi}(y) \bar{\alpha}(y) \mathbb{1}_C(y, y) \nu(dy) \\ = \iint h_{\pi}(y) \delta_y(dx) \bar{\alpha}(y) \mathbb{1}_C(x, y) \nu(dy) . \end{aligned} \quad (4.8)$$

Hence, summing (4.7) and (4.8) we obtain

$$\iint h_{\pi}(x) P(x, dy) \nu(dx) \mathbb{1}_C(x, y) = \iint h_{\pi}(y) P(y, dx) \mathbb{1}_C(x, y) \nu(dy) .$$

This proves that  $\pi$  is reversible with respect to  $P$ . ■

From Lemma 4.2, we obtain that  $\pi$  is an invariant probability for the Markov kernel  $P$ .

**Example 4.4 (Random walk Metropolis algorithm).** This is a particular case of the Metropolis-Hastings algorithm, where the proposal transition density is symmetric, i.e.  $q(x, y) = q(y, x)$ , for every  $(x, y) \in \mathbf{X} \times \mathbf{X}$ . Furthermore, assume that  $\mathbf{X} = \mathbb{R}^d$  and let  $\bar{q}$  be a symmetric density with respect to 0, i.e.  $\bar{q}(-y) = \bar{q}(y)$  for all  $y \in \mathbf{X}$ . Consider the transition density  $q$  defined by  $q(x, y) = \bar{q}(y - x)$ . This means that if the current state is  $X_k$ , an increment  $Z_{k+1}$  is drawn from  $\bar{q}$  and the candidate  $Y_{k+1} = X_k + Z_{k+1}$  is proposed.

The acceptance probability (4.3) for the random walk Metropolis algorithm is given by

$$\alpha(x, y) = 1 \wedge \frac{h_{\pi}(y)}{h_{\pi}(x)} . \quad (4.9)$$

If  $h_{\pi}(Y_{k+1}) \geq h_{\pi}(X_k)$ , then the move is accepted with probability one and if  $h_{\pi}(Y_{k+1}) < h_{\pi}(X_k)$ , then the move is accepted with a probability strictly less than one.

The choice of the incremental distribution is crucial for the efficiency of the algorithm. A classical choice for  $\bar{q}$  is the multivariate normal distribution with zero-mean and covariance matrix  $\Gamma$  to be suitably chosen.

**Example 4.5 (Independent Metropolis-Hastings sampler).** Another possibility is to set the transition density to be  $q(x, y) = \bar{q}(y)$ , where  $\bar{q}$  is a density on  $\mathbf{X}$ . In this case, the next candidate is drawn independently of the current state of the chain. This yields the so-called independent sampler, which is closely related to the accept-reject algorithm for random variable simulation.

The acceptance probability (4.3) is given by

$$\alpha(x, y) = 1 \wedge \frac{h_{\pi}(y) \bar{q}(x)}{h_{\pi}(x) \bar{q}(y)} . \quad (4.10)$$

Candidate steps with a low weight  $\bar{q}(Y_{k+1})/\pi(Y_{k+1})$  are rarely accepted, whereas candidates with a high weight are very likely to be accepted. Therefore the chain will remain at these states for several steps with a high probability, thus increasing the importance of these states within the constructed sample.

Assume for example that  $h$  is the standard Gaussian density and that  $q$  is the density of the Gaussian distribution with zero mean and variance  $\sigma^2$ , so that  $q(x) = h(y/\sigma)/\sigma$ . Assume that  $\sigma^2 > 1$  so that the values being proposed are sampled from a distribution with heavier tails than the objective distribution  $h$ . Then the acceptance probability is

$$\alpha(x, y) = \begin{cases} 1 & |y| \leq |x|, \\ \exp(-(y^2 - x^2)(1 - \sigma^{-2})/2) & |y| > |x|. \end{cases}$$

Thus the algorithm accepts all moves which decrease the norm of the current state but only some of those which increase it.

If  $\sigma^2 < 1$ , the values being proposed are sampled from a lighter tailed distribution than  $h$  and the acceptance probability becomes

$$\alpha(x, y) = \begin{cases} \exp(-(x^2 - y^2)(1 - \sigma^{-2})/2) & |y| \leq |x|, \\ 1 & |y| > |x|. \end{cases}$$

It is natural to inquire whether heavy-tailed or light-tailed proposal distributions should be preferred. This question will be partially answered in ??.

**Example 4.6 (Langevin diffusion).** More sophisticated proposals can be considered. Assuming that  $x \mapsto \log h_\pi(x)$  is everywhere differentiable. Consider the Langevin diffusion defined by the stochastic differential equation (SDE)

$$dX_t = \frac{1}{2} \nabla \log h_\pi(X_t) dt + dW_t,$$

where  $\nabla \log h_\pi$  denotes the gradient of  $\log h_\pi$ . Under appropriate conditions, the Langevin diffusion has a stationary distribution with density  $h_\pi$  and is reversible. Assume that the proposal state  $Y_{k+1}$  in the Metropolis-Hastings algorithm corresponds to the Euler discretization of the Langevin SDE for some step size  $h$ :

$$Y_{k+1} = X_k + \frac{\gamma}{2} \nabla \log h_\pi(X_k) + \sqrt{\gamma} Z_{k+1}, \quad Z_{k+1} \sim N(0, I).$$

Such algorithms are known as Langevin Metropolis-Hastings algorithms. The gradient can be approximated numerically via finite differences and does not require knowledge of the normalizing constant of the target distribution  $h_\pi$ .

### 4.1.2 Data augmentation

Throughout this section,  $(X, \mathcal{X})$  and  $(Y, \mathcal{Y})$  are Polish spaces equipped with their Borel  $\sigma$ -fields. Again, we wish to simulate from a probability measure  $\pi$  defined on  $(X, \mathcal{X})$  using a sequence  $\{X_k, k \in \mathbb{N}\}$  of  $X$ -valued random variables. Data augmentation algorithms consist in writing the target distribution  $\pi$  as the marginal of the distribution  $\pi^*$  on the product space  $(X \times Y, \mathcal{X} \otimes \mathcal{Y})$  defined by  $\pi^* = \pi \otimes R$  where  $R$  is a kernel on  $X \times \mathcal{Y}$ . By ??, there exists also a kernel  $S$  on  $Y \times \mathcal{X}$  and a probability measure  $\tilde{\pi}$  on  $(Y, \mathcal{Y})$  such that  $\pi^*(C) = \iint \mathbb{1}_C(x, y) \tilde{\pi}(dy) S(y, dx)$  for  $C \in \mathcal{X} \otimes \mathcal{Y}$ . In other words, if  $(X, Y)$  is a pair of random variables with distribution  $\pi^*$ , then  $R(x, \cdot)$  is the distribution of  $Y$  conditionally on  $X = x$  and  $S(y, \cdot)$  is the distribution of  $X$  conditionally on  $Y = y$ . The bivariate distribution  $\pi^*$  can then be expressed as follows

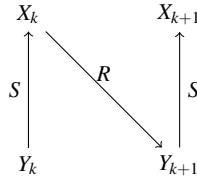
$$\pi^*(dxdy) = \pi(dx)R(x, dy) = S(y, dx)\tilde{\pi}(dy). \quad (4.11)$$



A data augmentation algorithm consists in running a Markov Chain  $\{(X_k, Y_k), k \in \mathbb{N}\}$  with invariant probability  $\pi^*$  and to use  $n^{-1} \sum_{k=0}^{n-1} f(X_k)$  as an approximation of  $\pi(f)$ . A significant difference between this general approach and a Metropolis-Hastings algorithm associated to the target distribution  $\pi$  is that  $\{X_k, k \in \mathbb{N}\}$  is no longer constrained to be a Markov chain. The transition from  $(X_k, Y_k)$  to  $(X_{k+1}, Y_{k+1})$  is decomposed into two successive steps:  $Y_{k+1}$  is first drawn given  $(X_k, Y_k)$  and then  $X_{k+1}$  is drawn given  $(X_k, Y_{k+1})$ . Intuitively,  $Y_{k+1}$  can be used as an auxiliary variable, which directs the moves of  $X_k$  toward interesting regions with respect to the target distribution.

When sampling from  $R$  and  $S$  is feasible, a classical choice consists in following the two successive steps: given  $(X_k, Y_k)$ ,

- (i) sample  $Y_{k+1}$  from  $R(X_k, \cdot)$ ,
- (ii) sample  $X_{k+1}$  from  $S(Y_{k+1}, \cdot)$ .



**Fig. 4.1** In this example, sampling from  $R$  and  $S$  is feasible.

It turns out that  $\{X_k, k \in \mathbb{N}\}$  is a Markov chain with Markov kernel  $RS$  and  $\pi$  is reversible with respect to  $RS$ .

**Lemma 4.7** *The distribution  $\pi$  is reversible with respect to the kernel  $RS$ .*

PROOF. By ??, we must prove that the measure  $\pi \otimes RS$  on  $X^2$  is symmetric. For  $A, B \in \mathcal{X}$ , applying (4.11), we have

$$\begin{aligned} \pi \otimes RS(A \times B) &= \int_{X \times Y} \pi(dx) R(x, dy) \mathbb{1}_A(x) S(y, B) = \int_{X \times Y} \mathbb{1}_A(x) S(y, B) \pi^*(dx dy) \\ &= \int_{X \times Y} \mathbb{1}_A(x) S(y, B) S(y, dx) \tilde{\pi}(dy) = \int_Y S(y, A) S(y, B) \tilde{\pi}(dy). \end{aligned}$$

This proves that  $\pi \otimes RS$  is symmetric. ■

Assume now that sampling from  $R$  or  $S$  is infeasible. In this case, we consider two instrumental kernels  $Q$  on  $(X \times Y) \times \mathcal{Y}$  and  $T$  on  $(X \times Y) \times \mathcal{X}$  which will be used to propose successive candidates for  $Y_{k+1}$  and  $X_{k+1}$ . For simplicity, assume that  $R(x, dy')$  and  $Q(x, y; dy')$  (resp.  $S(y', dx')$  and  $T(x, y'; dx')$ ) are dominated by the same measure and call  $r$  and  $q$  (resp.  $s$  and  $t$ ) the associated transition densities. We assume that  $r$  and  $s$  are known up to a normalizing constant. Define the Markov chain  $\{(X_k, Y_k), k \in \mathbb{N}\}$  as follows. Given  $(X_k, Y_k) = (x, y)$ ,

- (DA1) draw a candidate  $\tilde{Y}_{k+1}$  according to the distribution  $Q(x, y; \cdot)$  and accept  $Y_{k+1} = \tilde{Y}_{k+1}$  with probability  $\alpha(x, y, \tilde{Y}_{k+1})$  defined by

$$\alpha(x, y, y') = \frac{r(x, y') q(x, y'; y)}{r(x, y) q(x, y; y')} \wedge 1;$$

otherwise, set  $Y_{k+1} = Y_k$ ; the Markov kernel on  $X \times Y \times \mathcal{Y}$  associated to this transition is denoted by  $K_1$ ;

- (DA2) draw then a candidate  $\tilde{X}_{k+1}$  according to the distribution  $T(x, Y_{k+1}; \cdot)$  and accept  $X_{k+1} = \tilde{X}_{k+1}$  with probability  $\beta(x, Y_{k+1}, \tilde{X}_{k+1})$  defined by

$$\beta(x, y, x') = \frac{s(y, x') t(x', y; x)}{s(y, x) t(x, y; x')} \wedge 1;$$

otherwise, set  $X_{k+1} = X_k$ ; the Markov kernel on  $X \times Y \times \mathcal{X}$  associated to this transition is denoted by  $K_2$ .

For  $i = 1, 2$ , let  $K_i^*$  be the kernels associated to  $K_1$  and  $K_2$  as follows: for  $x \in X$ ,  $y \in Y$ ,  $A \in \mathcal{X}$  and  $B \in \mathcal{Y}$ ,

$$K_1^*(x, y; A \times B) = \mathbb{1}_A(x) K_1(x, y; B) . \quad (4.12)$$

$$K_2^*(x, y; A \times B) = \mathbb{1}_B(y) K_2(x, y; A) . \quad (4.13)$$

Then, the kernel of the chain  $\{(X_n, Y_n), n \in \mathbb{N}\}$  is  $K = K_1^* K_2^*$ . The process  $\{X_n, n \in \mathbb{N}\}$  is in general not a Markov chain since the distribution of  $X_{k+1}$  conditionally on  $(X_k, Y_k)$  depends on  $(X_k, Y_k)$  and on  $X_k$  only, except in some special cases. Obviously, this construction includes the previous one where sampling from  $R$  and  $S$  was feasible. Indeed, if  $Q(x, y; \cdot) = R(x, \cdot)$  and  $T(x, y; \cdot) = S(x, \cdot)$ , then the acceptance probabilities  $\alpha$  and  $\beta$  defined above simplify to one, the candidates are always accepted and we are back to the previous algorithm.

**Proposition 4.8** *The extended target distribution  $\pi^*$  is reversible with respect to the kernels  $K_1^*$  and  $K_2^*$  and invariant with respect to  $K$ .*

PROOF. The reversibility of  $\pi^*$  with respect to  $K_1^*$  and  $K_2^*$  implies its invariance and consequently its invariance with respect to the product  $K = K_1^* K_2^*$ . Let us prove the reversibility of  $\pi^*$  with respect to  $K_1^*$ . For each  $x \in X$ , the kernel  $K_1(x, \cdot; \cdot)$  on  $Y \times \mathcal{Y}$  is the kernel of a Metropolis-Hastings algorithm with target density  $r(x, \cdot)$ , proposal kernel density  $q(x, \cdot; \cdot)$  and acceptance probability  $\alpha(x, \cdot, \cdot)$ . By Proposition 4.3, this implies that the distribution  $R(x, \cdot)$  is reversible with respect to the kernel  $K_1(x, \cdot; \cdot)$ . Applying the definition (4.12) of  $K_1^*$  and  $\pi^* = \pi \otimes R$  yields, for  $A, C \in \mathcal{X}$  and  $B, D \in \mathcal{Y}$ ,

$$\begin{aligned} \pi^* \otimes K_1^*(A \times B \times C \times D) &= \iint_{A \times B} \pi(\mathrm{d}x \mathrm{d}y) K_1^*(x, y; C \times D) \\ &= \iint_{A \times B} \pi(\mathrm{d}x) R(x, \mathrm{d}y) \mathbb{1}_C(x) K_1(x, y, D) \\ &= \int_{A \cap C} \pi(\mathrm{d}x) [R(x, \cdot) \otimes K_1(x, \cdot; \cdot)](B \times D) . \end{aligned}$$

We have seen that for each  $x \in X$ , the measure  $R(x, \cdot) \otimes K_1(x, \cdot; \cdot)$  is symmetric, thus  $\pi^* \otimes K^*$  is also symmetric. The reversibility of  $\pi^*$  with respect to  $K_2^*$  is proved similarly. ■

**Example 4.9 (The slice sampler).** Set  $X = \mathbb{R}^d$  and  $\mathcal{X} = \mathcal{B}(X)$ . Let  $\mu$  be a  $\sigma$ -finite measure on  $(X, \mathcal{X})$  and let  $h$  be the density with respect to  $\mu$  of the target distribution. We assume that for all  $x \in X$ ,

$$h(x) = C \prod_{i=0}^k f_i(x) ,$$

where  $C$  is a constant (which is not necessarily known) and  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}_+$  are nonnegative measurable functions. For  $y = (y_1, \dots, y_k) \in \mathbb{R}_+^k$ , define

$$L(y) = \left\{ x \in \mathbb{R}^d : f_i(x) \geq y_i, i = 1, \dots, k \right\} .$$

The  $f_0$ -slice-sampler algorithm proceeds as follows:

- given  $X_n$ , draw independently a  $k$ -tuple  $Y_{n+1} = (Y_{n+1,1}, \dots, Y_{n+1,k})$  of independent random variables such that  $Y_{n+1,i} \sim \text{Unif}(0, f_i(X_n))$ ,  $i = 1, \dots, k$ .
- sample  $X_{n+1}$  from the distribution with density proportional to  $f_0 \mathbb{1}_{L(Y_{n+1})}$ .

Set  $Y = \mathbb{R}_+^k$  and for  $(x, y) \in X \times Y$ ,

$$h^*(x, y) = C f_0(x) \mathbb{1}_{L(y)}(x) = h(x) \prod_{i=1}^k \frac{\mathbb{1}_{[0, f_i(x)]}(y_i)}{f_i(x)} .$$

Let  $\pi^*$  be the probability measure with density  $h^*$  with respect to Lebesgue's measure on  $X \times Y$ . Then  $\int_Y h^*(x, y) dy = h(x)$  i.e.  $\pi$  is the first marginal of  $\pi^*$ . Let  $R$  be the kernel on  $X \times \mathcal{Y}$  with kernel density  $r$  defined by

$$r(x, y) = \frac{h^*(x, y)}{h(x)} \mathbb{1}_{\{h(x) > 0\}} .$$

Then  $\pi^* = \pi \otimes R$ . Define the distribution  $\tilde{\pi} = \pi R$ , its density  $\tilde{h}(y) = \int_X h^*(u, y) du$  and the kernel  $S$  on  $Y \times \mathcal{X}$  with density  $s$  by

$$s(y, x) = \frac{h^*(x, y)}{\tilde{h}(y)} \mathbb{1}_{\{\tilde{h}(y) > 0\}} .$$

If  $(X, Y)$  is a vector with distribution  $\pi^*$ , then  $S(y, \cdot)$  is the conditional distribution of  $X$  given  $Y = y$  and the Markov kernel of the chain  $\{X_n, n \in \mathbb{N}\}$  is  $RS$  and Lemma 4.7 can be applied to prove that  $\pi$  is reversible, hence invariant, with respect to  $RS$ .

### 4.1.3 Two-stage Gibbs sampler

The Gibbs sampler is a simple method which decomposes a complex multidimensional distribution into a collection of smaller dimensional ones. Let  $(X, \mathcal{X})$  and  $(Y, \mathcal{Y})$  be complete separable metric spaces endowed with their Borel  $\sigma$ -fields. To construct the Markov chain  $\{(X_n, Y_n), n \in \mathbb{N}\}$  with  $\pi^*$  as an invariant probability, we proceed exactly as in data-augmentation algorithms. Assume that  $\pi^*$  may be written as

$$\pi^*(dxdy) = \pi(dx)R(x, dy) = \tilde{\pi}(dy)S(y, dx) \quad (4.14)$$

where  $\pi$  and  $\tilde{\pi}$  are probability measures on  $X$  and  $Y$  respectively and  $R$  and  $S$  are kernels on  $X \times \mathcal{Y}$  and  $Y \times \mathcal{X}$  respectively.

#### The deterministic updating (two-stage) Gibbs (DUGS) sampler

When sampling from  $R$  and  $S$  is feasible, the DUGS sampler proceeds as follows: given  $(X_k, Y_k)$ ,

- (DUGS1) sample  $Y_{k+1}$  from  $R(X_k, \cdot)$ ,
- (DUGS2) sample  $X_{k+1}$  from  $S(Y_{k+1}, \cdot)$ .

For both the Data Augmentation algorithms and the two-stage Gibbs sampler we consider a distribution  $\pi^*$  on the product space  $X \times Y$ . In the former case, the distribution of interest is a marginal distribution of  $\pi^*$  and in the latter case the target distribution is  $\pi^*$  itself.

We may associate to each update (DUGS1)-(DUGS2) of the algorithm a transition kernel on  $(X \times Y) \times (\mathcal{X} \otimes \mathcal{Y})$  defined for  $(x, y) \in X \times Y$  and  $A \times B \in \mathcal{X} \otimes \mathcal{Y}$  by

$$R^*(x, y; A \times B) = \mathbb{1}_A(x)R(x, B) , \quad (4.15)$$

$$S^*(x, y; A \times B) = \mathbb{1}_B(y)S(y, A) . \quad (4.16)$$

The transition kernel of the DUGS is then given by

$$P_{\text{DUGS}} = R^* S^* . \quad (4.17)$$

Note that for  $A \times B \in \mathcal{X} \otimes \mathcal{Y}$ ,

$$\begin{aligned}
P_{\text{DUGS}}(x, y; A \times B) &= \iint_{X \times Y} R^*(x, y; dx' dy') S^*(x', y'; A \times B) \\
&= \iint_{X \times Y} R(x, dy') \mathbb{1}_B(y') S(y', A) \\
&= \int_B R(x, dy') S(y', A) = R \otimes S(x, B \times A). \tag{4.18}
\end{aligned}$$

As a consequence of Proposition 4.8, we obtain the invariance of  $\pi^*$ .

**Lemma 4.10** *The distribution  $\pi^*$  is reversible with respect to the kernels  $R^*$  and  $S^*$  and invariant with respect to  $P_{\text{DUGS}}$ .*

### The Random Scan Gibbs sampler (RSGS)

At each iteration, the RSGS algorithm consists in updating one component chosen at random. It proceeds as follows: given  $(X_k, Y_k)$ ,

(RSGS1) sample a Bernoulli random variable  $B_{k+1}$  with probability of success  $1/2$ .

(RSGS2) If  $B_{k+1} = 0$ , then sample  $Y_{k+1}$  from  $R(X_k, \cdot)$  else sample  $X_{k+1}$  from  $S(Y_{k+1}, \cdot)$ .

The transition kernel of the RSGS algorithm can be written

$$P_{\text{RSGS}} = \frac{1}{2} R^* + \frac{1}{2} S^*. \tag{4.19}$$

Lemma 4.10 implies that  $P_{\text{RSGS}}$  is reversible with respect to  $\pi^*$  and therefore  $\pi^*$  is invariant for  $P_{\text{RSGS}}$ .

If sampling from  $R$  or  $S$  is infeasible, the Gibbs transitions can be replaced by a Metropolis-Hastings algorithm on each component as in the case of the DUGS algorithm. The algorithm is then called the Two-Stage Metropolis-within-Gibbs algorithm.

**Example 4.11 (Scalar Normal-Inverse Gamma).** In a statistical problem one may be presented with a set of independent observations  $\mathbf{y} = \{y_1, \dots, y_n\}$ , assumed to be normally distributed, but with unknown mean  $\mu$  and variance  $\tau^{-1}$  ( $\tau$  is often referred to as the precision). The Bayesian approach to this problem is to assume that  $\mu$  and  $\tau$  are themselves random variables, with a given prior distribution. For example, we might assume that

$$\mu \sim N(\theta_0, \phi_0^{-1}), \quad \tau \sim \Gamma(a_0, b_0), \tag{4.20}$$

i.e.  $\mu$  is normally distributed with mean  $\theta_0$  and variance  $\phi_0^{-1}$  and  $\tau$  has a Gamma distribution with parameters  $a_0$  and  $b_0$ .

The parameters  $\theta_0, \phi_0, a_0$  and  $b_0$  are assumed to be known. The posterior density  $h$  of  $(\mu, \tau)$  defined as the conditional density given the observations, is then given, using the Bayes formula, by

$$h(u, t) \propto \exp(-\phi_0(u - \theta_0)^2/2) \exp\left(-t \sum_{i=1}^n (y_i - u)^2/2\right) t^{a_0-1+n/2} \exp(-b_0 t).$$

Conditioning on the observations introduces a dependence between  $\mu$  and  $\tau$ . Nevertheless, the conditional laws of  $\mu$  given  $\tau$  and  $\tau$  given  $\mu$  have a simple form. Write  $\bar{y} = n^{-1} \sum_{i=1}^n y_i$  and  $S^2 = n^{-1} \sum_{i=1}^n (y_i - \bar{y})^2$ ,

$$\begin{aligned}
\theta_n(t) &= (\phi_0 \theta_0 + nt \bar{y}) / (\phi_0 + nt), \quad \phi_n(t) = \phi_0 + nt, \\
a_n &= a_0 + n/2, \quad b_n(u) = b_0 + nS^2/2 + n(\bar{y} - u)^2/2.
\end{aligned}$$

Then,

$$\mathcal{L}(\mu|\tau) = N(\theta_n(\tau), \phi_n^{-1}(\tau)), \quad \mathcal{L}(\tau|\mu) = \Gamma(a_n, b_n(\mu)).$$

The Gibbs sampler provides a simple approach to define a Markov chain whose invariant probability has the density  $h$ . First we simulate  $\mu_0$  and  $\tau_0$  independently with distribution as in (4.20). At the  $k$ -th stage,

given  $(\mu_{k-1}, \tau_{k-1})$ , we first simulate  $N_k \sim N(0, 1)$  and  $G_k \sim \Gamma(a_n, 1)$  and we set

$$\begin{aligned}\mu_k &= \theta_n(\tau_{k-1}) + \phi_n^{-1/2}(\tau_{k-1})N_k \\ \tau_k^{-1} &= b_n(\mu_k)G_k.\end{aligned}$$

In the simple case where  $\theta_0 = 0$  and  $\phi_0 = 0$  which corresponds to a flat prior for  $\mu$  (an improper distribution with a constant density on  $\mathbb{R}$ ), the above equation can be rewritten as

$$\begin{aligned}\mu_k &= \bar{y} + (n\tau_{k-1})^{-1/2}N_k \\ \tau_k^{-1} &= (b_0 + nS^2/2 + n(\bar{y} - \mu_k)^2/2)G_k.\end{aligned}$$

Thus,  $\{\tau_k^{-1}, k \in \mathbb{N}\}$  and  $\{(\mu_k - \bar{y})^2, k \in \mathbb{N}\}$  are Markov chains which follow the random coefficient autoregressions

$$\begin{aligned}\tau_k^{-1} &= \frac{N_k^2 G_k}{2} \tau_{k-1}^{-1} + \left(b_0 + \frac{nS^2}{2}\right) G_k, \\ (\mu_k - \bar{y})^2 &= \frac{N_k^2 G_{k-1}}{2} (\mu_{k-1} - \bar{y})^2 + \left(b_0 + \frac{nS^2}{2}\right) G_{k-1}.\end{aligned}$$

#### 4.1.4 Hit-and-run algorithm

Let  $K$  be a bounded subset of  $\mathbb{R}^d$  with non-empty interior. Let  $\rho : K \rightarrow [0, \infty)$  be a (not necessarily normalized) density, i.e. a non-negative Lebesgue-integrable function. We define the probability measure  $\pi_\rho$  with density  $\rho$  by

$$\pi_\rho(A) = \frac{\int_A \rho(x) dx}{\int_K \rho(x) dx} \quad (4.21)$$

for all measurable sets  $A \subset K$ . For example, if  $\rho(x) \equiv 1$  then  $\pi$  is simply the uniform distribution on  $K$ . The hit-and-run Markov kernel, presented below, can be used to sample approximately from  $\pi_\rho$ . The hit-and-run algorithm consists of two steps. Starting from  $x \in K$ , we first choose a random direction  $\theta \in S_{d-1}$ , the unit sphere in  $\mathbb{R}^d$  according to a uniform distribution on  $S_{d-1}$ . We then choose the next state of the Markov chain with respect to the density  $\rho$  restricted to the chord determined by the current state  $x$  and the direction  $\theta \in S_{d-1}$ : for any function  $f \in \mathbb{F}_+(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ ,

$$H_\theta f(x) = \frac{1}{\ell_\rho(x, \theta)} \int_{s=-\infty}^{\infty} \mathbb{1}_K(x + s\theta) f(x + s\theta) \rho(x + s\theta) ds, \quad (4.22)$$

where  $\ell_\rho(x, \theta)$  is the normalizing constant defined as

$$\ell_\rho(x, \theta) = \int_{-\infty}^{\infty} \mathbb{1}_K(x + s\theta) \rho(x + s\theta) ds. \quad (4.23)$$

The Markov kernel  $H$  that corresponds to the hit-and-run algorithm is therefore defined by, for all  $f \in \mathbb{F}_+(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  and  $x \in \mathbb{R}^d$ ,

$$Hf(x) = \int_{S_{d-1}} H_\theta f(x) \sigma_{d-1}(d\theta), \quad (4.24)$$

where  $\sigma_{d-1}$  is the uniform distribution on  $S_{d-1}$ .

**Lemma 4.12** *For all  $\theta \in S_{d-1}$ , the Markov kernel  $H_\theta$  is reversible with respect to  $\pi_\rho$  defined in (4.21). Furthermore,  $H$  is also reversible with respect to  $\pi_\rho$ .*

PROOF. Let  $c = \int_K \rho(x) dx$  and  $A, B \in \mathcal{B}(K)$ . By elementary computations, we have

$$\begin{aligned}
\int_A H_\theta(x, B) \pi_\rho(dx) &= \int_A \int_{-\infty}^{\infty} \frac{\mathbb{1}_B(x+s\theta) \rho(x+s\theta) ds}{\ell_\rho(x, \theta)} \frac{\rho(x) dx}{c} \\
&= \int_{\mathbb{R}^d} \int_{-\infty}^{\infty} \frac{\mathbb{1}_A(x) \mathbb{1}_B(x+s\theta) \rho(x) ds dx}{\ell_\rho(x, \theta)} \frac{1}{c} \\
&= \int_{\mathbb{R}^d} \int_{-\infty}^{\infty} \frac{\mathbb{1}_A(y-s\theta) \mathbb{1}_B(y) \rho(y) ds dy}{\ell_\rho(y-s\theta, \theta)} \frac{1}{c} \\
&= \int_B \int_{-\infty}^{\infty} \frac{\mathbb{1}_A(y-s\theta) \rho(y-s\theta)}{\ell_\rho(y-s\theta, \theta)} ds \pi_\rho(dy) \\
&= \int_B \int_{-\infty}^{\infty} \frac{\mathbb{1}_A(y-s\theta) \rho(y-s\theta)}{\ell_\rho(y, \theta)} ds \pi_\rho(dy) \\
&= \int_B \int_{-\infty}^{\infty} \frac{\mathbb{1}_A(y+t\theta) \rho(y+t\theta)}{\ell_\rho(y, \theta)} dt \pi_\rho(dy) \\
&= \int_B H_\theta(x, A) \pi_\rho(dx) .
\end{aligned}$$

The reversibility of  $H_\theta$  is proved. The reversibility of  $H_\rho$  follows from the reversibility of  $H_\theta$ : for any  $A, B \in \mathcal{B}(\mathbb{R}^d)$ , we have

$$\begin{aligned}
\int_A H(x, B) \pi_\rho(dx) &= \int_{S_{d-1}} \int_A H_\theta(x, B) \pi_\rho(dx) \sigma_{d-1}(d\theta) \\
&= \int_{S_{d-1}} \int_B H_\theta(x, A) \pi_\rho(dx) \sigma_{d-1}(d\theta) \\
&= \int_B H(x, A) \pi_\rho(dx) .
\end{aligned}$$

■

### 4.1.5 Gibbs sampling

When the distribution of interest is multivariate, it may be the case that for each particular variable, its conditional distribution given all remaining variables has a simple form. In this case, a natural algorithm is the *Gibbs sampler*, which is now described. Its name somehow inappropriately stems from its use for the simulation of Gibbs Markov random fields by ?.

Assume that  $\mathbf{X} = \mathbf{X}_1 \times \cdots \times \mathbf{X}_m$  is a product space equipped with the product  $\sigma$ -algebra  $\mathcal{X} = \mathcal{X}_1 \otimes \cdots \otimes \mathcal{X}_m$ . Let  $\lambda_k$ ,  $k \in \{1, \dots, m\}$ , be  $\sigma$ -finite measures on  $(\mathbf{X}_k, \mathcal{X}_k)$ , and let  $\lambda = \lambda_1 \otimes \lambda_2 \otimes \cdots \otimes \lambda_m$  be the product measure. Suppose we are given a joint distribution with probability density function  $\pi$  with respect to the product measure  $\lambda$ . For simplicity, we assume that  $\pi$  is everywhere positive. An element  $x \in \mathbf{X}$  may be decomposed into  $m$  components  $x = (x^{[1]}, \dots, x^{[m]})$ , where  $x^{[k]} \in \mathbf{X}_k$ . If  $k$  is an index in  $\{1, \dots, m\}$ , we shall denote by  $x^{[k]}$  the  $k$ th component of  $x$  and by  $x^{[-k]} = \{x^{[l]}\}_{l \neq k}$  the collection of remaining components. We further denote by  $\pi_k(\cdot \mid x^{[-k]})$  the conditional probability density function, defined as

$$\pi_k(x^{[k]} \mid x^{[-k]}) = \frac{\pi(x^{[1]}, x^{[2]}, \dots, x^{[m]})}{\int \pi(x^{[1]}, x^{[2]}, \dots, x^{[m]}) \lambda_k(dx^{[k]})} .$$

We assume that sampling from this conditional distribution is feasible (for  $k = 1, \dots, m$ ). Note that  $x^{[k]}$  is not necessarily scalar but may be itself vector-valued.

The *deterministic scan Gibbs sampler* is an MCMC algorithm which, starting from an initial arbitrary state  $X_0$ , updates the current state  $X_i = (X_i^{[1]}, \dots, X_i^{[m]})$  to a new state  $X^{[i+1]}$  as follows.

For  $k = 1, 2, \dots, m$ : Simulate  $X_{i+1}^{[k]}$  from  $\pi_k(\cdot \mid X_{i+1}^{[1]}, \dots, X_{i+1}^{[k-1]}, X_i^{[k+1]}, \dots, X_i^{[m]})$ .

In other words, in the  $k$ th round of the cycle needed to simulate  $X_{i+1}$ , the  $k$ th component is updated by simulation from its conditional distribution given all other components, which remain fixed. This new value then supersedes the old one and is used in the subsequent simulation steps. A complete cycle of  $m$  condi-

tional simulations is usually referred to as a *sweep* of the algorithm. We denote by  $K_k$  the corresponding Markov kernel:

$$K_k(x, A_1 \times \cdots \times A_m) = K_k(x^{[1]}, \dots, x^{[m]}; A) \\ = \int \cdots \int \prod_{j \neq k} \delta_{x^{[j]}}(dx'^{[j]}) \pi_k(x'^{[k]} | x^{[-k]}) \mathbb{1}_A(x') \lambda_k(dx'^{[k]}) . \quad (4.25)$$

**Proposition 4.13 (Reversibility of individual Gibbs steps)** *Each of the  $m$  individual kernels  $K_k$ ,  $k \in \{1, \dots, m\}$  is reversible with respect to  $\pi$  and thus admits  $\pi$  as a stationary probability density function.*

PROOF. See Exercise 8.26. ■

Each step corresponds to a very special type of Metropolis-Hastings move where the acceptance probability is equal to 1, due to choice of  $\pi_k$  as the proposal distribution. However, Proposition 4.13 does not suffice to establish the convergence of the Gibbs sampler. Only the combination of the  $m$  moves in the complete cycle has a chance of producing a chain with the ability to visit the whole space  $X$  from any starting point.

**Example 4.14 (Scalar normal-inverse gamma)** *In a statistical problem one may be presented with a set of independent observations  $\mathbf{y} := \{y_1, \dots, y_n\}$ , which we assume to be normally distributed, but with unknown mean  $\mu$  and variance  $\tau^{-1}$  ( $\tau$  is often referred to as the precision). The Bayesian approach to this problem is to assume that  $\mu$  and  $\tau$  are themselves random variables, with a given prior distribution. For example, we might assume that*

$$\mu \sim N(\mu_0, \tau_0^{-1}), \quad \tau \sim \Gamma(a_0, b_0),$$

i.e.,  $\mu$  is normal with mean  $\mu_0$  and variance  $\tau_0^{-1}$  and  $\tau$  has gamma distribution with parameters  $a_0$  and  $b_0$ .<sup>1</sup>

The parameters  $\mu_0$ ,  $\tau_0$ ,  $a_0$  and  $b_0$  are assumed to be known. Then the prior density for  $(\mu, \tau)$  is given by

$$\pi(\mu, \tau) \propto \exp\{-\tau_0(\mu - \mu_0)^2/2\} \tau^{a_0-1} \exp\{-b_0\tau\}.$$

The posterior density for  $(\mu, \tau)$  defined as the conditional density given the observations, is then given, using the Bayes formula, by

$$\pi(\mu, \tau | \mathbf{y}) \propto \pi(\mu, \tau) f(\mathbf{y} | \mu, \tau) \propto \exp(-\tau_0(\mu - \mu_0)^2/2) \\ \times \exp\left(-\tau \sum_{i=1}^n (y_i - \mu)^2/2\right) \tau^{a_0-1+n/2} \exp(-b_0\tau).$$

Conditioning with respect to the observations has introduced a dependence between  $\mu$  and  $\tau$ . Nevertheless, the full conditional distributions still have a simple form

$$\pi(\mu | \mathbf{y}, \tau) = N(m_n(\tau), t_n^{-1}(\tau)), \\ \pi(\tau | \mathbf{y}, \mu) = \Gamma(a_n, b_n(\mu)),$$

where, denoting  $\bar{y} = n^{-1} \sum_{i=1}^n y_i$

$$m_n(\tau) := (\tau_0 \mu_0 + n \tau \bar{y}) / (\tau_0 + n \tau), \quad t_n(\tau) := \tau_0 + n \tau, \\ a_n := a_0 + n/2, \quad b_n(\mu) := b_0 + 1/2 \sum_{i=1}^n (y_i - \mu)^2.$$

<sup>1</sup>  $Z$  has an inverse gamma distribution if  $1/Z$  has a gamma distribution; general properties can be found, for example, in (?), Section 8.5)

The Gibbs sampler provides a particularly simple approach to sample from  $\pi(\mu, \tau \mid \mathbf{y})$ . We set

$$\mathbf{X} = \mathbf{X}_1 \times \mathbf{X}_2 = \mathbb{R} \times [0, \infty) .$$

We wish to simulate  $X = (\mu, \tau)$  with density  $\pi(\mu, \tau \mid \mathbf{y})$ . First we simulate  $X_0$ , say from the product form density  $\pi(\mu, \tau)$ . At the  $t$ -th stage, given  $X_{t-1} = (\mu_{t-1}, \tau_{t-1})$ , we first simulate  $N_t \sim N(0, 1)$  and  $G_t \sim \Gamma(a_n, 1)$  and we put

$$\begin{aligned} \mu_t &= m_n(\tau_{t-1}) + t_n^{-1/2}(\tau_{t-1})N_t \\ \tau_t &= b_n(\mu_t)G_t , \end{aligned}$$

and  $X_t = (\mu_t, \tau_t)$ . In the simple case where  $\mu_0 = 0$  and  $\tau_0 = 0$ , which corresponds to a flat prior for  $\mu$  (an improper distribution with a “constant” density on  $\mathbb{R}$ ), the above equation can be rewritten as

$$\begin{aligned} \mu_t &= \bar{y} + (n\tau_{t-1})^{-1/2}N_t \\ \tau_t &= \left( b_0 + 1/2 \sum_{i=1}^n (y_i - \mu_t)^2 \right) G_t . \end{aligned}$$

By alternating and iterating these two updates, we see that  $\tau_t^{-1}$  is a Markov chain that iterates following a random coefficient autoregression (see ??),

$$\tau_t^{-1} = A_t \tau_{t-1}^{-1} + B_t , \quad (4.26)$$

where  $A_t = N_t^2 / (2G_t)$ , and  $B_t = (b_0 + nS^2/2) / G_t$ . It is easily checked that  $(\mu_t - \bar{y})^2$  is also a random coefficient autoregressive process. The Gibbs sampler can be implemented in **R** using the script provided for this example.

## 4.2 Ordering the asymptotic variances

Define  $L^2(\pi) = \{f : X \rightarrow \mathbb{R} : f \text{ is } \mathcal{X}/\mathcal{B}(\mathbb{R})\text{-measurable}\}$  and  $L_0^2(\pi) = \{f \in L^2(\pi) : \pi(f) = 0\}$ . Moreover, for all  $f, g \in L^2(\pi)$ , write

$$\langle f, g \rangle = \pi(fg) .$$

**Definition 4.15.** Let  $P_0$  and  $P_1$  be Markov transition kernels on  $(X, \mathcal{X})$  with invariant probability  $\pi$ . We say that  $P_1$  dominates  $P_0$  on the off-diagonal, written  $P_0 \leq P_1$ , if for all  $A \in \mathcal{X}$ , and  $\pi$ -a.e. all  $x$  in  $X$ ,

$$P_0(x, A \setminus \{x\}) \leq P_1(x, A \setminus \{x\}) .$$

**Definition 4.16.** Let  $P_0$  and  $P_1$  be Markov transition kernels on  $(X, \mathcal{X})$  with invariant probability  $\pi$ . We say that  $P_1$  dominates  $P_0$  in the covariance ordering, written  $P_0 \preceq P_1$ , if for all  $f \in L^2(\pi)$ ,

$$\langle f, P_1 f \rangle \leq \langle f, P_0 f \rangle .$$

**Lemma 4.17** Let  $P_0$  and  $P_1$  be Markov transition kernels on  $(X, \mathcal{X})$  with invariant probability  $\pi$ . Assume moreover that  $P_0 \leq P_1$ . Then,  $P_0 \preceq P_1$ .

PROOF. For all  $x \in X$  and  $A \in \mathcal{X}$ , define



$$P(x, A) = \delta_x(A) + P_1(x, A) - P_0(x, A) .$$

$P$  is a nonnegative kernel since for all  $x \in \mathbf{X}$  and  $A \in \mathcal{X}$ ,

$$\begin{aligned} P(x, \{x\}) &= 1 + P_1(x, \{x\}) - P_0(x, \{x\}) \geq 0 , \\ P(x, A \setminus \{x\}) &= P_1(x, A \setminus \{x\}) - P_0(x, A \setminus \{x\}) \geq 0 . \end{aligned}$$

Combining with  $P(x, \mathbf{X}) = 1$ , this implies that  $P$  is a Markov kernel. We now show that for all  $f \in \mathbf{L}^2(\pi)$

$$\langle f, P_0 f \rangle - \langle f, P_1 f \rangle = \iint \pi(\mathrm{d}x) P(x, \mathrm{d}y) (f(x) - f(y))^2 / 2$$

Indeed,

$$\begin{aligned} \langle f, P_0 f \rangle - \langle f, P_1 f \rangle &= \iint \pi(\mathrm{d}x) f(x) (P_0(x, \mathrm{d}y) - P_1(x, \mathrm{d}y)) f(y) \\ &= \iint \pi(\mathrm{d}x) f(x) (\delta_x(\mathrm{d}y) - P(x, \mathrm{d}y)) f(y) \\ &= \int \pi(\mathrm{d}x) f^2(x) - \iint \pi(\mathrm{d}x) P(x, \mathrm{d}y) f(x) f(y) \\ &= \iint \pi(\mathrm{d}x) P(x, \mathrm{d}y) \left[ \frac{f^2(x) - f^2(y)}{2} + f(x) f(y) \right] \end{aligned}$$

where the last inequality follows from the fact that  $P$  is clearly  $\pi$ -invariant. Finally,

$$\langle f, P_0 f \rangle - \langle f, P_1 f \rangle = \iint \pi(\mathrm{d}x) P(x, \mathrm{d}y) (f(x) - f(y))^2 / 2 \geq 0 .$$

And thus,  $P_0 \preceq P_1$ . ■

**Proposition 4.18** *Let  $P_0$  and  $P_1$  be Markov kernels on  $\mathbf{X} \times \mathcal{X}$  and  $\pi \in \mathbb{M}_1(\mathcal{X})$ . Assume that  $\pi$  is reversible with respect to  $P_0$  and  $P_1$  and that for all  $i \in \{0, 1\}$ ,  $\pi(|f|^2) + 2 \lim_{n \rightarrow \infty} \sum_{k=1}^n \pi(|f| P_i^k |f|) < \infty$ . Assume that  $P_0 \preceq P_1$ . Then, for any  $f \in \mathbf{L}_0^2(\pi)$ ,*

$$v_1(f, P_1) \leq v_0(f, P_0) ,$$

where for  $i \in \{0, 1\}$

$$v_i(f, P_i) = \pi(f^2) + 2 \lim_{n \rightarrow \infty} \sum_{k=1}^n \pi(f P_i^k f) .$$

PROOF. For all  $\alpha \in (0, 1)$ , denote  $P_\alpha = (1 - \alpha)P_0 + \alpha P_1$ . For  $\lambda \in (0, 1)$ , define

$$w_\lambda(\alpha) = \sum_{k=0}^{\infty} \lambda^k \langle f, P_\alpha^k f \rangle ,$$

We first show that for all  $\alpha \in (0, 1)$ ,

$$\frac{\mathrm{d}w_\lambda(\alpha)}{\mathrm{d}\alpha} = \sum_{k=0}^{\infty} \lambda^k \sum_{i=1}^k \langle f, P_\alpha^{i-1} (P_1 - P_0) P_\alpha^{k-i} f \rangle . \quad (4.27)$$

Indeed, for all  $1 \leq \ell \leq k$  and all  $\alpha_1, \dots, \alpha_k$ ,

$$\begin{aligned} \langle f, P_{\alpha_1} \dots P_{\alpha_k} f \rangle &= (1 - \alpha_\ell) \langle f, P_{\alpha_1} \dots P_{\alpha_{\ell-1}} P_0 P_{\alpha_{\ell+1}} \dots P_{\alpha_k} f \rangle \\ &\quad + \alpha_\ell \langle f, P_{\alpha_1} \dots P_{\alpha_{\ell-1}} P_1 P_{\alpha_{\ell+1}} \dots P_{\alpha_k} f \rangle , \end{aligned}$$

so that

$$\frac{\partial}{\partial \alpha_\ell} \langle f, P_{\alpha_1} \dots P_{\alpha_k} f \rangle = \langle f, P_{\alpha_1} \dots P_{\alpha_{\ell-1}} (P_1 - P_0) P_{\alpha_{\ell+1}} \dots P_{\alpha_k} f \rangle ,$$

and thus, we obtain by differentiating  $\alpha \mapsto w_\lambda(\alpha)$ ,

$$\frac{\mathrm{d}w_\lambda(\alpha)}{\mathrm{d}\alpha} = \sum_{k=0}^{\infty} \lambda^k \sum_{i=1}^k \langle f, P_\alpha^{i-1} (P_1 - P_0) P_\alpha^{k-i} f \rangle .$$

This shows (4.27). Using now that  $\pi$  is reversible for the kernel  $P_\alpha$ ,

$$\begin{aligned}
\frac{dw_\lambda(\alpha)}{d\alpha} &= \sum_{i=1}^{\infty} \sum_{k \geq i}^{\infty} \lambda^k \left\langle P_\alpha^{i-1} f, (P_1 - P_0) P_\alpha^{k-i} f \right\rangle \\
&= \lambda \left\langle \sum_{\ell=0}^{\infty} \lambda^\ell P_\alpha^\ell f, (P_1 - P_0) \sum_{\ell=0}^{\infty} \lambda^\ell P_\alpha^\ell f \right\rangle \leq 0,
\end{aligned}$$

which completes the proof. ■

# Chapter 5

## Ergodic theory for Markov chains

### Contents

<b>5.1</b>	<b>Dynamical systems</b> .....	<b>80</b>
5.1.1	Definitions .....	80
5.1.2	Invariant events .....	81
<b>5.2</b>	<b>Markov chains ergodicity</b> .....	<b>84</b>
<b>5.3</b>	<b>Exercises</b> .....	<b>89</b>
<b>5.4</b>	<b>Bibliographical notes</b> .....	<b>93</b>

This chapter is concerned with the asymptotic behaviour of sample averages of stationary ergodic Markov chains. For this purpose, it is convenient to link the Markov chain to a certain dynamical system. The Law of Large Numbers for Markov chains is then obtained as a consequence of the classical Birkhoff theorem. It turns out that under appropriate assumptions, this approach still holds true for functions that actually depend on the whole trajectory, such as  $n^{-1} \sum_{k=0}^{n-1} f(\{X_{k+\ell}, \ell \in \mathbb{N}\})$  or  $n^{-1} \sum_{k=0}^{n-1} f(\{X_{k-\ell}, \ell \in \mathbb{N}\})$ . A key result of this chapter is Theorem 5.21 which shows that the existence of a unique invariant probability measure implies the ergodicity of the associated dynamical system, which in turn allows to apply the Birkhoff ergodic theorem. Still, the price to pay for using the dynamical system theory is the stationarity assumption. Typically in this chapter, the Law of Large Numbers will be proved  $\mathbb{P}_\pi$  – a.s. (where  $\pi$  is the unique invariant probability measure for  $P$ ) and will be extended to other initial distributions. Sufficient conditions are given in this chapter but a more thorough treatment for other initial distributions requires notions that will be introduced in later chapters.

## 5.1 Dynamical systems

We first briefly introduce some basic definitions and properties of dynamical systems that will be useful when applying them to Markov chains.

### 5.1.1 Definitions

**Definition 5.1 (Dynamical system).** Let  $(\Omega, \mathcal{B}, \mathbb{P})$  be a probability space.

- A measurable map  $T$  from  $(\Omega, \mathcal{B})$  to  $(\Omega, \mathcal{B})$  is a measure-preserving transformation if for all  $A \in \mathcal{B}$ ,

$$\mathbb{P}(T^{-1}(A)) = \mathbb{P}(A) .$$

The probability  $\mathbb{P}$  is then said to be invariant under the transformation  $T$  and  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  is said to be a dynamical system.

- The application  $T$  is said to be an invertible measure-preserving transformation if it is measure-preserving, invertible and its inverse  $T^{-1}$  is measurable.

If the transformation  $T$  is measure preserving and invertible, then  $T^{-1}$  is also measure-preserving since, for all  $A \in \mathcal{B}$ ,

$$\mathbb{P}((T^{-1})^{-1}(A)) = \mathbb{P}(T(A)) = \mathbb{P}(T^{-1}\{T(A)\}) = \mathbb{P}(A) .$$

Note also that if  $T$  is measure-preserving, then for all integer  $n \in \mathbb{N}$  and  $A \in \mathcal{B}$ ,

$$\mathbb{P}(T^{-n}(A)) = \mathbb{P}(A) .$$

Let  $(X, \mathcal{X})$  be a measurable space. Denote by  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  the associated canonical space and by  $\{X_n, n \in \mathbb{N}\}$  the coordinate process. The shift operator  $\theta$  (see ??) is defined, for  $\omega = (\omega_k)_{k \in \mathbb{N}} \in X^{\mathbb{N}}$ , by

$$\theta(\omega_0, \omega_1, \dots) = (\omega_1, \omega_2, \dots) .$$

Note that  $X_k \circ \theta = X_{k+1}$  for all  $k \geq 0$ . By ??,  $\theta$  is a measurable map from  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  to  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$ , but it is not invertible. Recall that  $\{X_n, n \in \mathbb{N}\}$  is stationary if the distribution of  $(X_k, \dots, X_{k+n})$  is independent of  $k$  for all  $n \in \mathbb{N}$ . The next Lemma shows the connection between the stationarity of the coordinate process and the invariance of  $\mathbb{P}$  under the shift operator  $\theta$ .

**Lemma 5.2** *A probability measure  $\mathbb{P}$  on  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  is invariant under the shift operator  $\theta$  if and only if the coordinate process is stationary under  $\mathbb{P}$ .*

PROOF. It suffices to note (??) that  $\mathbb{P}$  is measure preserving if and only if for all  $n \geq 1$  and all  $f \in \mathbb{F}_b(X^n, \mathcal{X}^{\otimes n})$ ,  $\mathbb{E}[f(X_0, \dots, X_{n-1})] = \mathbb{E}[f(X_1, \dots, X_n)]$ . ■

**Example 5.3 (One-sided Markov shift).** Let  $P$  be a kernel on  $(X, \mathcal{X})$  which admits an invariant probability  $\pi$  on  $(X, \mathcal{X})$ . By ??, there exists a unique probability measure  $\mathbb{P}_\pi$  on the canonical space such that the coordinate process is a Markov chain with kernel  $P$  and initial distribution  $\pi$ . Then, by ??, the canonical chain is a stationary process. Lemma 5.2 then shows that  $\mathbb{P}_\pi$  is invariant under  $\theta$ , i.e.  $\mathbb{P}_\pi \circ \theta^{-1} = \mathbb{P}_\pi$ .

### 5.1.2 Invariant events

**Definition 5.4 (Invariant random variable, invariant event).** Let  $T$  be a measurable map from  $(\Omega, \mathcal{B})$  to  $(\Omega, \mathcal{B})$ .

- A  $\mathbb{R}$ -valued random variable  $Y$  on  $(\Omega, \mathcal{B})$  is invariant for  $T$  if  $Y \circ T = Y$ .
- An event  $A$  is invariant for  $T$  if  $A = T^{-1}(A)$  or equivalently if its indicator function  $\mathbb{1}_A$  is invariant for  $T$ .

**Proposition 5.5** *Let  $T$  be a measurable map from  $(\Omega, \mathcal{B})$  to  $(\Omega, \mathcal{B})$ .*

- (i) *The collection  $\mathcal{I}$  of invariant sets for  $T$  is a sub- $\sigma$ -field of  $\mathcal{B}$ .*
- (ii) *Let  $Y$  be a  $\mathbb{R}$ -valued random variable.  $Y$  is invariant if and only if  $Y$  is  $\mathcal{I}$ -measurable.*

PROOF. The proof of (i) is elementary and omitted. Consider now (ii). If  $Y \circ T = Y$ , then for all  $B \in \mathcal{B}(\mathbb{R})$ ,

$$T^{-1}(Y^{-1}(B)) = (Y \circ T)^{-1}(B) = Y^{-1}(B).$$

Thus  $Y^{-1}(B) \in \mathcal{I}$  and  $Y$  is  $\mathcal{I}$ -measurable.

Conversely, if  $Y$  is  $\mathcal{I}$ -measurable, define  $A_{k,n} = \{\frac{k}{n} \leq Y < \frac{k+1}{n}\} \in \mathcal{I}$ ,  $n \geq 1$ ,  $k \in \mathbb{Z}$ . Then, with the convention  $\infty \times 0 = 0$ ,  $Y$  is the pointwise limit of the sequence  $\{Y_n, n \in \mathbb{N}^*\}$  defined by

$$Y_n = \sum_{k \in \mathbb{Z}} \frac{k}{n} \mathbb{1}_{A_{k,n}} + \infty \mathbb{1}_{\{Y = +\infty\}} - \infty \mathbb{1}_{\{Y = -\infty\}}.$$

Since  $Y$  is  $\mathcal{I}$ -measurable, the sets  $A_{k,n}$ ,  $\{Y = -\infty\}$  and  $\{Y = +\infty\}$  belong to  $\mathcal{I}$  hence the functions  $Y_n$  are invariant for all  $n$  and

$$Y \circ T = \left( \lim_{n \rightarrow \infty} Y_n \right) \circ T = \lim_{n \rightarrow \infty} (Y_n \circ T) = \lim_{n \rightarrow \infty} Y_n = Y.$$

■

The most important examples of invariant random variables which we will be considered in the sequel are defined as limits.

**Lemma 5.6** *Let  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  be the canonical space,  $\{X_n, n \in \mathbb{N}\}$  the coordinate process and  $\theta$  the shift operator. Then  $\mathcal{I} \subset \bigcap_{k \geq 0} \sigma(X_\ell, \ell \geq k)$ . Moreover, for any  $f \in \mathbb{F}(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ ,  $\limsup_{n \rightarrow \infty} f(X_n)$ ,  $\liminf_{n \rightarrow \infty} f(X_n)$ ,  $\limsup_{n \rightarrow \infty} n^{-1}(f(X_0) + \dots + f(X_{n-1}))$  and  $\liminf_{n \rightarrow \infty} n^{-1}(f(X_0) + \dots + f(X_{n-1}))$  are invariant random variables.*

PROOF. Set  $\mathcal{G}_k = \sigma(X_\ell, \ell \geq k)$  and  $\mathcal{G}_\infty = \bigcap_{k \geq 0} \mathcal{G}_k$ . Let  $A$  be an invariant set. Then,  $A \in \mathcal{G}_0 = \mathcal{X}^{\otimes \mathbb{N}}$ . Since we have the implication: if  $A \in \mathcal{G}_k$ , then  $A = \theta^{-1}(A) \in \mathcal{G}_{k+1}$ , we obtain by induction that  $A \in \mathcal{G}_k$  for all  $k$  and thus  $A \in \bigcap_{k \geq 0} \mathcal{G}_k$ .

The remaining statements of the lemma are straightforward.  $\blacksquare$

Let now  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be a dynamical system, that is,  $T$  is measure preserving for  $\mathbb{P}$ . A  $\mathbb{R}$ -valued random variable  $Y$  defined on  $\Omega$  is said to be  $\mathbb{P}$ -a.s. invariant (for  $T$ ) if  $Y \circ T = Y$ ,  $\mathbb{P}$ -a.s. Similarly, an event  $A \in \mathcal{B}$  is  $\mathbb{P}$ -a.s. invariant (for  $T$ ) if its indicator function  $\mathbb{1}_A$  is  $\mathbb{P}$ -a.s. invariant.

**Lemma 5.7** *If  $Y$  is  $\mathbb{P}$ -a.s. invariant, then there exists an invariant random variable  $Z$ , such that  $Y = Z$   $\mathbb{P}$ -a.s. In particular, if  $A \in \mathcal{B}$  is  $\mathbb{P}$ -a.s. invariant, there exists  $B \in \mathcal{I}$  such that  $\mathbb{1}_A = \mathbb{1}_B$   $\mathbb{P}$ -a.s.*

PROOF. The random variable  $Z = \limsup_{n \rightarrow \infty} Y \circ T^n$  is invariant. Since  $Y$  is  $\mathbb{P}$ -a.s. invariant,  $Y = Y \circ T$   $\mathbb{P}$ -a.s., hence  $Y = Y \circ T^n$   $\mathbb{P}$ -a.s. for all  $n \geq 1$ . This yields  $Y = Z$   $\mathbb{P}$ -a.s. If  $Y = \mathbb{1}_A$ , then there exists an invariant random variable  $Z$  such that  $\mathbb{1}_A = Z$   $\mathbb{P}$ -a.s.. The set  $B = \{Z = 1\}$  is therefore invariant and  $\mathbb{1}_A = \mathbb{1}_B$   $\mathbb{P}$ -a.s..  $\blacksquare$

It is easy to check that the family of  $\mathbb{P}$ -a.s. invariant sets for  $T$  is a  $\sigma$ -algebra  $\mathcal{I}_{\mathbb{P}}$ . Lemma 5.7 shows that  $\mathcal{I}_{\mathbb{P}}$  is the  $\mathbb{P}$ -completion of the invariant  $\sigma$ -algebra  $\mathcal{I}$  (the  $\sigma$ -algebra generated by  $\mathcal{I}$  and the family of sets which are  $\mathbb{P}$ -negligible).

Denote by  $T^n$  the transformation  $T$  iterated  $n$ -times and by convention, we let  $T^0$  be the identity function. The behavior of time averages is given by the following fundamental result.

**Theorem 5.8 (Birkhoff's ergodic theorem).** *Let  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be a dynamical system and  $Y$  be a random variable such that  $\mathbb{E}[|Y|] < \infty$ . Then,*

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \mathbb{E}[Y | \mathcal{I}] \quad \mathbb{P}\text{-a.s.} \quad (5.1)$$

Moreover, the convergence also holds in  $L^1(\mathbb{P})$ .

The proof is based on the following lemma.

**Lemma 5.9** *Let  $Z$  be a random variable such that  $\mathbb{E}[|Z|] < \infty$ . If  $\mathbb{E}[Z | \mathcal{I}] > 0$   $\mathbb{P}$ -a.s., then,*

$$\liminf_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Z \circ T^k \geq 0 \quad \mathbb{P}\text{-a.s.}$$

PROOF. For all  $n \in \mathbb{N}^*$ , write  $S_n = \sum_{k=0}^{n-1} Z \circ T^k$ . Note that for all  $n \geq 1$ ,  $\mathbb{E}[|S_n|] \leq n\mathbb{E}[|Z|] < \infty$ . Denote  $L_n = \inf\{S_k : 1 \leq k \leq n\}$  and  $A = \{\inf_{n \in \mathbb{N}^*} L_n = -\infty\}$ . Since  $|Z| < \infty$   $\mathbb{P}$ -a.s.,  $\{\inf_{n \geq 1} S_n = -\infty\} = \{\inf_{n \geq 1} S_n \circ T = -\infty\}$   $\mathbb{P}$ -a.s., the set  $A$  is  $\mathbb{P}$ -a.s. invariant. Since  $L_{n-1} \geq L_n$ ,

$$\begin{aligned} L_n &= Z + \inf\{S_k - Z : 1 \leq k \leq n\} \\ &= Z + \inf(0, L_{n-1} \circ T) \geq Z + \inf(0, L_n \circ T). \end{aligned} \quad (5.2)$$

Since  $\mathbb{E}[|S_k|] < \infty$  for all  $k \in \mathbb{N}$ , for all  $n \geq 1$ ,  $\mathbb{E}[|L_n|] \leq \sum_{k=0}^{n-1} \mathbb{E}[|S_k|] < \infty$ . and (5.2) implies that  $Z \leq L_n + (L_n \circ T)^- = L_n + L_n^- \circ T$   $\mathbb{P}$ -a.s. Then, using  $\mathbb{1}_A = \mathbb{1}_A \circ T$   $\mathbb{P}$ -a.s., we get

$$\begin{aligned} \mathbb{E}[\mathbb{1}_A Z] &\leq \mathbb{E}[\mathbb{1}_A L_n] + \mathbb{E}[\mathbb{1}_A L_n^- \circ T] = \mathbb{E}[\mathbb{1}_A L_n] + \mathbb{E}[\mathbb{1}_A \circ T L_n^- \circ T] \\ &\leq \mathbb{E}[\mathbb{1}_A L_n] + \mathbb{E}[\mathbb{1}_A L_n^-] = \mathbb{E}[\mathbb{1}_A L_n^+]. \end{aligned} \quad (5.3)$$

Since  $L_n^+ \leq Z^+$  with  $\mathbb{E}[Z^+] < \infty$  and  $\lim_{n \rightarrow \infty} \mathbb{1}_A L_n^+ = 0$   $\mathbb{P}$ -a.s., Lebesgue's dominated convergence theorem shows that  $\mathbb{E}[\lim_{n \rightarrow \infty} \mathbb{1}_A L_n^+] = 0$ . Therefore, since  $0 \leq \mathbb{E}[\mathbb{1}_A \mathbb{E}[Z | \mathcal{I}]] = \mathbb{E}[\mathbb{1}_A Z]$ , we finally get using (5.3)

$$\mathbb{E}[\mathbb{1}_A \mathbb{E}[Z | \mathcal{I}]] = \mathbb{E}[\mathbb{1}_A Z] \leq \mathbb{E}\left[\lim_{n \rightarrow \infty} \mathbb{1}_A L_n^+\right] = 0.$$

By assumption,  $\mathbb{E}[Z | \mathcal{I}] > 0$   $\mathbb{P}$ -a.s., the previous inequality shows  $\mathbb{P}(A) = 0$ . We conclude that  $\liminf_{n \rightarrow \infty} n^{-1} S_n \geq 0$   $\mathbb{P}$ -a.s.  $\blacksquare$

PROOF. [of Theorem 5.8] Let  $\varepsilon > 0$  and set  $Z = Y - \mathbb{E}[Y | \mathcal{I}] + \varepsilon$ . Note that  $\mathbb{E}[|Z|] \leq 2\mathbb{E}[|Y|] + \varepsilon$  showing that  $\mathbb{E}[Z | \mathcal{I}]$

is well-defined and, by construction,  $\mathbb{E}[Z|\mathcal{J}] > 0$ . Using that  $\mathbb{E}[Y|\mathcal{J}]$  being  $\mathcal{J}$ -measurable, it is invariant according to Proposition 5.5 and Lemma 5.9 implies that

$$\liminf_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k \geq \mathbb{E}[Y|\mathcal{J}] - \varepsilon \quad \mathbb{P} - \text{a.s.}$$

Replacing  $Y$  by  $-Y$ , we finally obtain

$$-\varepsilon + \limsup_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k \leq \mathbb{E}[Y|\mathcal{J}] \leq \liminf_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k + \varepsilon \quad \mathbb{P} - \text{a.s.}$$

This shows (5.1) since  $\varepsilon > 0$  is arbitrary.

We now turn to the  $L^1(\mathbb{P})$  convergence. Denote by  $M_n(Y) = n^{-1} \sum_{k=0}^{n-1} Y \circ T^k$ . If  $Y$  bounded, then Lebesgue's dominated convergence theorem shows that the convergence in (5.1) also holds in  $L^1(\mathbb{P})$ . For a bounded random variable  $\bar{Y}$ , consider the decomposition

$$|M_n(Y) - \mathbb{E}[Y|\mathcal{J}]| \leq |M_n(Y) - M_n(\bar{Y})| + |M_n(\bar{Y}) - \mathbb{E}[\bar{Y}|\mathcal{J}]| + \mathbb{E}[|\bar{Y} - Y||\mathcal{J}] .$$

Let us denote by  $\|U\|_1 = \mathbb{E}[|U|]$ . Note that  $\|\mathbb{E}[|\bar{Y} - Y||\mathcal{J}]\|_1 \leq \|\bar{Y} - Y\|_1$  and

$$\|M_n(Y) - M_n(\bar{Y})\|_1 \leq n^{-1} \sum_{k=0}^{n-1} \|(Y - \bar{Y}) \circ T^k\|_1 = \|Y - \bar{Y}\|_1 ,$$

where we have used that the transformation  $T$  is measure preserving and thus  $\|(Y - \bar{Y}) \circ T^k\|_1 = \|Y - \bar{Y}\|_1$  for all  $k \in \mathbb{N}$ . Therefore,

$$\limsup_{n \rightarrow \infty} \|M_n(Y) - \mathbb{E}[Y|\mathcal{J}]\|_1 \leq 2\|\bar{Y} - Y\|_1 .$$

The proof is complete since bounded random variables are dense in  $L^1(\mathbb{P})$ . ■

The most interesting case of application of Theorem 5.8 is when the  $\sigma$ -field  $\mathcal{J}$  is trivial, in which case the conditional expectation  $\mathbb{E}[Y|\mathcal{J}]$  can be replaced by  $\mathbb{E}[Y]$  in (5.1).

**Definition 5.10 (Ergodic dynamical system).** A dynamical system  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  is ergodic if the invariant  $\sigma$ -field  $\mathcal{J}$  is trivial for  $\mathbb{P}$ , i.e. for all  $A \in \mathcal{J}$ ,  $\mathbb{P}(A) \in \{0, 1\}$ .

**Corollary 5.11** *Let  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be an ergodic dynamical system and  $Y$  be a  $\mathbb{R}$ -valued random variable such that  $\mathbb{E}[|Y|] < \infty$ . Then,*

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \mathbb{E}[Y] \quad \mathbb{P} - \text{a.s.} \quad (5.4)$$

### 5.1.2.1 Dynamical systems associated to one-sided and two-sided sequences

In the context of Markov chains on a measurable space  $(X, \mathcal{X})$ , the dynamical systems will be associated to either one-sided sequences  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$ , or two-sided sequences  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}})$ . Denote by  $\bar{\theta} : X^{\mathbb{Z}} \rightarrow X^{\mathbb{Z}}$  the shift operator on  $X^{\mathbb{Z}}$ : for any two-sided sequence  $\omega = (\omega_n)_{n \in \mathbb{Z}} \in X^{\mathbb{Z}}$ ,

$$[\bar{\theta}(\omega)]_n = \omega_{n+1} , \quad \text{for all } n \in \mathbb{Z} . \quad (5.5)$$

Moreover, set  $\bar{\theta}_1 = \bar{\theta}$  and for all  $n > 1$ ,  $\bar{\theta}_n = \bar{\theta}_{n-1} \circ \bar{\theta}$ . Let  $\bar{\mathbb{P}}$  be a probability measure on  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}})$  and denote by  $\bar{\mathbb{E}}$  the associated expectation operator. Let  $\Pi$  be the measurable map from  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}})$  to  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  defined by

$$\Pi(\omega) = (\omega_n)_{n \in \mathbb{N}} , \quad \text{for all } \omega = (\omega_n)_{n \in \mathbb{Z}} \in X^{\mathbb{Z}} . \quad (5.6)$$

We denote by  $\mathbb{P} = \bar{\mathbb{P}} \circ \Pi^{-1}$  the probability induced on  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  by the probability  $\bar{\mathbb{P}}$  and the map  $\Pi$ . Let  $\theta$  be the shift operator on  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  defined in ?? and note that  $\theta \circ \Pi = \Pi \circ \bar{\theta}$ .

**Lemma 5.12** *If  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}}, \bar{\mathbb{P}}, \bar{\theta})$  is a dynamical system, then  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}, \theta)$  is also a dynamical system.*

PROOF. By assumption,  $\bar{\mathbb{P}} \circ \bar{\theta}^{-1} = \bar{\mathbb{P}}$ . Combining with  $\theta \circ \Pi = \Pi \circ \bar{\theta}$ , this implies for all  $A \in \mathcal{X}^{\otimes \mathbb{N}}$ ,

$$\begin{aligned} \mathbb{P} \circ \theta^{-1}(A) &= (\bar{\mathbb{P}} \circ \Pi^{-1}) \circ \theta^{-1}(A) = \bar{\mathbb{P}} \circ (\theta \circ \Pi)^{-1}(A) = \bar{\mathbb{P}} \circ (\Pi \circ \bar{\theta})^{-1}(A) \\ &= (\bar{\mathbb{P}} \circ \bar{\theta}^{-1}) \circ \Pi^{-1}(A) = \bar{\mathbb{P}} \circ \Pi^{-1}(A) = \mathbb{P}(A). \end{aligned}$$

■

**Proposition 5.13** *Assume that  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}}, \bar{\mathbb{P}}, \bar{\theta})$  is a dynamical system. If the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}, \theta)$  is ergodic, then  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}}, \bar{\mathbb{P}}, \bar{\theta})$  is ergodic.*

PROOF.

Let  $A$  be an invariant set for the dynamical system  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}}, \bar{\mathbb{P}}, \bar{\theta})$ , that is  $\mathbb{1}_A = \mathbb{1}_A \circ \bar{\theta}$ . We will show that  $\bar{\mathbb{P}}(A) = 0$  or  $1$ .

Note first that  $\mathcal{X}^{\otimes \mathbb{Z}} = \sigma(\mathcal{F}_{-k}^+, k \in \mathbb{N})$  where  $\mathcal{F}_{\ell}^+ = \sigma(X_i, \ell \leq i < \infty)$  and  $\{X_i, i \in \mathbb{Z}\}$  are the coordinate process on  $X^{\mathbb{Z}}$ . This allows to apply the approximation ?? showing that for all  $\varepsilon > 0$ , there exists  $k_{\varepsilon} \in \mathbb{N}^*$  and a  $\mathcal{F}_{-k_{\varepsilon}}^+$ -measurable random variable  $Z_{\varepsilon}$  such that  $\bar{\mathbb{E}}[|Z_{\varepsilon}|] < \infty$  and  $\bar{\mathbb{E}}[|\mathbb{1}_A - Z_{\varepsilon}|] \leq \varepsilon$ . Set  $Y_{\varepsilon} = Z_{\varepsilon} \circ \bar{\theta}_{k_{\varepsilon}}$ . By construction,  $Y_{\varepsilon}$  is  $\mathcal{F}_0^+$ -measurable. Using that  $A$  is an invariant set, we obtain

$$\bar{\mathbb{E}}[|\mathbb{1}_A - Y_{\varepsilon}|] = \bar{\mathbb{E}}[|\mathbb{1}_A \circ \bar{\theta}_{k_{\varepsilon}} - Z_{\varepsilon} \circ \bar{\theta}_{k_{\varepsilon}}|] = \bar{\mathbb{E}}[|\mathbb{1}_A - Z_{\varepsilon}|] \leq \varepsilon.$$

Since  $\varepsilon$  is arbitrary, there exists a  $\mathcal{F}_0^+$ -measurable random variable  $Y$  satisfying  $\bar{\mathbb{E}}[|Y|] < \infty$  and  $\mathbb{1}_A = Y$ ,  $\bar{\mathbb{P}}$ -a.s. Since  $1 = \bar{\mathbb{P}}(\mathbb{1}_A = Y) \leq \bar{\mathbb{P}}(Y \in \{0, 1\}) \leq 1$  there exists  $B \in \mathcal{F}_0^+$  such that

$$\mathbb{1}_B = Y = \mathbb{1}_A, \quad \bar{\mathbb{P}}\text{-a.s.} \quad (5.7)$$

Eq. (5.7) and the invariance of  $A$  then shows that

$$\bar{\mathbb{P}}(\mathbb{1}_B \circ \bar{\theta} = \mathbb{1}_A \circ \bar{\theta} = \mathbb{1}_A = \mathbb{1}_B) = 1.$$

Now, note that  $\mathcal{F}_0^+ = \sigma(\Pi)$ , the  $\sigma$ -algebra generated by  $\Pi$ , where the canonical projection  $\Pi : X^{\mathbb{Z}} \rightarrow \Omega$  is defined in (5.6). Then, since  $B \in \mathcal{F}_0^+$ , there exists  $C \in \mathcal{F}^+$  such that  $B = \Pi^{-1}(C)$  and thus,

$$\begin{aligned} 1 &= \bar{\mathbb{P}}(\mathbb{1}_B = \mathbb{1}_B \circ \bar{\theta}) \\ &= \bar{\mathbb{P}}(\mathbb{1}_C \circ \Pi = \mathbb{1}_C \circ \Pi \circ \bar{\theta}) \\ &\stackrel{(i)}{=} \bar{\mathbb{P}}(\mathbb{1}_C \circ \Pi = \mathbb{1}_C \circ \theta \circ \Pi) = \bar{\mathbb{P}} \circ \Pi^{-1}(\mathbb{1}_C = \mathbb{1}_C \circ \theta) \stackrel{(ii)}{=} \mathbb{P}(\mathbb{1}_C = \mathbb{1}_C \circ \theta), \end{aligned}$$

where  $\stackrel{(i)}{=}$  follows from  $\Pi \circ \bar{\theta} = \theta \circ \Pi$  and  $\stackrel{(ii)}{=}$  from  $\mathbb{P} = \bar{\mathbb{P}} \circ \Pi^{-1}$ . The dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}, \theta)$  being ergodic, it implies that  $\mathbb{P}(C) = 0$  or  $1$  which concludes the proof since

$$\mathbb{P}(C) = \bar{\mathbb{P}} \circ \Pi^{-1}(C) = \bar{\mathbb{P}}(B) = \bar{\mathbb{P}}(A).$$

■

Proposition 5.13 allows to study only the ergodicity of dynamical systems on one-sided sequences (instead of two-sided sequences) and then to use Birkhoff's ergodic theorem either to functions depending on the future  $n^{-1} \sum_{k=0}^{n-1} f(\{X_{k+\ell}, \ell \in \mathbb{N}\})$  or even on the whole past  $n^{-1} \sum_{k=0}^{n-1} f(\{X_{k-\ell}, \ell \in \mathbb{N}\})$ . From now on, we only consider the ergodicity of dynamical systems associated to one-sided sequences.

## 5.2 Markov chains ergodicity

We specialize the results of the previous section in the context of Markov chains. Here and subsequently, we consider a Markov kernel  $P$  on a measurable space  $(X, \mathcal{X})$  and the coordinate process  $\{X_k, k \in \mathbb{N}\}$  on the



canonical space  $(\Omega, \mathcal{F}) = (\mathcal{X}^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$ , endowed with the family of probability measures  $\mathbb{P}_\xi, \xi \in \mathbb{M}_1(\mathcal{X})$  under which the coordinate process is a Markov chain with kernel  $P$  and initial distribution  $\xi$ .

As a consequence of Birkhoff's ergodic theorem and of Corollary 5.11, we obtain the ergodic theorem for Markov chains.

**Theorem 5.14.** *Let  $P$  be a Markov kernel on  $\mathcal{X} \times \mathcal{X}$ . Assume that  $P$  admits an invariant probability measure  $\pi$  and that the associated dynamical system  $(\mathcal{X}^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic. Then, for all random variables  $Y \in \mathcal{L}^1(\mathbb{P}_\pi)$ ,*

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_\pi[Y] \quad \mathbb{P}_\pi - \text{a.s.}$$

Moreover, the convergence also holds in  $\mathcal{L}^1(\mathbb{P}_\pi)$ .

The condition  $Y \in \mathcal{L}^1(\mathbb{P}_\pi)$  may be relaxed to  $\mathbb{E}_\pi[Y^+] < \infty$  as shown by Exercise 5.33.

**Definition 5.15 (Absorbing set).** A set  $B \in \mathcal{X}$  is called absorbing if  $P(x, B) = 1$  for all  $x \in B$ .

This definition subsumes that the empty set is absorbing. Of course the interesting absorbing sets are non-empty.

**Proposition 5.16** *Let  $P$  be a Markov kernel on  $\mathcal{X} \times \mathcal{X}$  admitting an invariant probability measure  $\pi$ . If  $B \in \mathcal{X}$  is an absorbing set, then  $\pi_B = \pi(B \cap \cdot)$  is an invariant finite measure. Moreover, if the invariant probability measure is unique, then  $\pi(B) \in \{0, 1\}$ .*

PROOF. Let  $B$  be an absorbing set. Using that  $\pi_B \leq \pi$ ,  $\pi P = \pi$  and  $B$  is absorbing, we get that for all  $C \in \mathcal{X}$ ,

$$\pi_B P(C) = \pi_B P(C \cap B) + \pi_B P(C \cap B^c) \leq \pi P(C \cap B) + \pi_B P(B^c) = \pi(C \cap B) = \pi_B(C).$$

Replacing  $C$  by  $C^c$  and noting that  $\pi_B P(\mathcal{X}) = \pi_B(\mathcal{X}) < \infty$  show that  $\pi_B$  is an invariant finite measure. To complete the proof, assume that  $P$  has a unique invariant probability measure. If  $\pi(B) > 0$  then,  $\pi_B/\pi(B)$  is an invariant probability measure and is therefore equal to  $\pi$ . Since  $\pi_B(B^c) = 0$ , we get  $\pi(B^c) = 0$ . Thus,  $\pi(B) \in \{0, 1\}$ . ■

We now relate harmonic functions (defined in ??) with invariant random variables for the shift transformation  $\theta$ .

**Proposition 5.17** *Let  $P$  be a Markov kernel on  $\mathcal{X} \times \mathcal{X}$ .*

(i) *Let  $Y$  be a bounded invariant random variable for the shift transformation  $\theta$ . Then the function  $h_Y : x \mapsto h_Y(x) = \mathbb{E}_x[Y]$  is a bounded harmonic function.*

(ii) *Let  $h$  be a bounded harmonic function and define  $Y = \limsup_{n \rightarrow \infty} h(X_n)$ . Then  $Y$  is an invariant random variable for  $\theta$  and for any  $\xi \in \mathbb{M}_1(\mathcal{X})$ , the sequence  $\{h(X_n), n \in \mathbb{N}\}$  converges to  $Y$   $\mathbb{P}_\xi$  - a.s. and in  $\mathcal{L}^1(\mathbb{P}_\xi)$ . Moreover,  $h(x) = \mathbb{E}_x[Y]$  for all  $x \in \mathcal{X}$ .*

(iii) *Let  $\pi$  be an invariant probability measure and  $Y \in \mathcal{L}^1(\mathbb{P}_\pi)$  be an invariant random variable for  $\theta$ . Then,  $\mathbb{E}_x[|Y|] < \infty$   $\pi$  - a.e., the function  $x \mapsto \mathbb{E}_x[Y]$  is  $\pi$ -integrable and  $Y = \mathbb{E}_{X_0}[Y]$   $\mathbb{P}_\pi$  - a.s.*

PROOF.

- (i) Assume that  $Y : \mathbf{X}^{\mathbb{N}} \rightarrow \mathbb{R}$  is a bounded invariant random variable, i.e.  $Y \circ \theta = Y$ . By the Markov property, for any  $x \in \mathbf{X}$ ,

$$Ph_Y(x) = \mathbb{E}_x[h_Y(X_1)] = \mathbb{E}_x[\mathbb{E}_{X_1}[Y]] = \mathbb{E}_x[Y \circ \theta_1] = \mathbb{E}_x[Y] = h_Y(x),$$

showing that  $h_Y$  is harmonic.

- (ii) Let  $h$  be a bounded harmonic function:  $Ph(x) = h(x)$  for all  $x \in \mathbf{X}$ . Then,  $\{(h(X_n), \mathcal{F}_n), n \in \mathbb{N}\}$  is a bounded  $\mathbb{P}_\xi$ -martingale, for any initial distribution  $\xi \in \mathbb{M}_1(\mathcal{X})$ . By Doob's martingale convergence theorem, the sequence  $\{h(X_n), n \in \mathbb{N}\}$  converges  $\mathbb{P}_\xi$ -a.s. and in  $L^1(\mathbb{P}_\xi)$  to a limit. Hence, we get

$$Y = \lim_{n \rightarrow \infty} h(X_n) \quad \mathbb{P}_\xi - \text{a.s. and } \mathbb{E}_\xi[Y] = \lim_{n \rightarrow \infty} \mathbb{E}_\xi[h(X_n)]. \quad (5.8)$$

The function  $h$  being harmonic, we have  $h(x) = P^n h(x) = \mathbb{E}_x[h(X_n)]$  for all  $x \in \mathbf{X}$  and  $n \in \mathbb{N}$ . Applying (5.8) with  $\xi = \delta_x$

$$\mathbb{E}_x[Y] = \lim_{n \rightarrow \infty} \mathbb{E}_x[h(X_n)] = h(x).$$

- (iii) Since  $Y \in L^1(\mathbb{P}_\pi)$ ,  $\mathbb{E}_\pi[|Y|] = \int_{\mathbf{X}} \pi(dx) \mathbb{E}_x[|Y|]$  showing that  $\mathbb{E}_x[|Y|] < \infty$   $\pi$ -a.e. and that the function  $x \mapsto \mathbb{E}_x[Y]$  is integrable with respect to  $\pi$ . By the Markov property and the invariance of  $Y$ , we get

$$\mathbb{E}_{X_k}[Y] = \mathbb{E}_\pi[Y \circ \theta_k | \mathcal{F}_k] = \mathbb{E}_\pi[Y | \mathcal{F}_k] \quad \mathbb{P}_\pi - \text{a.s.}$$

Therefore,  $\{(\mathbb{E}_{X_k}[Y], \mathcal{F}_k), k \in \mathbb{N}\}$  is a uniformly integrable  $\mathbb{P}_\pi$ -martingale. By ??,

$$\lim_{k \rightarrow \infty} \mathbb{E}_{X_k}[Y] = \lim_{k \rightarrow \infty} \mathbb{E}_\pi[Y | \mathcal{F}_k] = \mathbb{E}_\pi[Y | \mathcal{F}] = Y \quad \mathbb{P}_\pi - \text{a.s.} \quad (5.9)$$

and in  $L^1(\mathbb{P}_\pi)$ . Moreover, applying successively that the translation operator  $\theta$  is measure preserving for  $\mathbb{P}_\pi$  and  $Y = Y \circ \theta_k$ , we obtain for any  $k \in \mathbb{N}$ ,

$$\begin{aligned} \mathbb{E}_\pi[|Y - \mathbb{E}_{X_0}[Y]|] &= \mathbb{E}_\pi[|Y - \mathbb{E}_{X_0}[Y] \circ \theta_k|] \\ &= \mathbb{E}_\pi[|Y \circ \theta_k - \mathbb{E}_{X_k}[Y]|] = \mathbb{E}_\pi[|Y - \mathbb{E}_{X_k}[Y]|]. \end{aligned}$$

Taking the limit as  $k$  goes to infinity, (5.9) yields

$$\mathbb{E}_\pi[|Y - \mathbb{E}_{X_0}[Y]|] = \lim_{k \rightarrow \infty} \mathbb{E}_\pi[|Y - \mathbb{E}_{X_k}[Y]|] = 0.$$

■

**Remark 5.18.** Proposition 5.17 shows that the map  $Y \mapsto h_Y$ , where  $h_Y(x) = \mathbb{E}_x[Y]$ ,  $x \in \mathbf{X}$  defines a one-to-one correspondence between the bounded harmonic functions and the bounded invariant random variables. If  $Y$  is a bounded invariant random variable, then  $h_Y : x \mapsto h_Y(x) = \mathbb{E}_x[Y]$  is a bounded harmonic function. If  $h$  is a bounded harmonic function, then  $h(x) = \mathbb{E}_x[Y]$  where  $Y = \limsup_{n \rightarrow \infty} h(X_n)$  is an invariant random variable (hence  $h = h_Y$ ). ▲

**Corollary 5.19** *Let  $P$  be a Markov kernel on  $\mathbf{X} \times \mathcal{X}$ . The following statements are equivalent.*

- (i) *The bounded harmonic functions are constant.*
- (ii) *The invariant  $\sigma$ -field  $\mathcal{I}$  is trivial up to an equivalence, i.e. for all  $A \in \mathcal{I}$ , we get for all  $\xi \in \mathbb{M}_1(\mathcal{X})$ ,  $\mathbb{P}_\xi(A) = 0$  or  $\mathbb{P}_\xi(A) = 1$ .*

PROOF. (i)  $\Rightarrow$  (ii): Let  $A$  be an invariant set. Then  $h_A : x \mapsto h_A(x) = \mathbb{P}_x(A)$  is a harmonic function by Proposition 5.17 which is constant under (i), i.e.  $h_A(x) = c$  for all  $x \in \mathbf{X}$ . By the Markov property, we get that, for all  $\xi \in \mathbb{M}_1(\mathcal{X})$ ,  $h_A(X_n) = \mathbb{E}_{X_n}[\mathbb{1}_A] = \mathbb{E}_\xi[\mathbb{1}_A \circ \theta_n | \mathcal{F}_n] \quad \mathbb{P}_\xi - \text{a.s.}$  Since  $A \in \mathcal{I}$ , it holds that  $\mathbb{E}_\xi[\mathbb{1}_A \circ \theta_n | \mathcal{F}_n] = \mathbb{E}_\xi[\mathbb{1}_A | \mathcal{F}_n] \quad \mathbb{P}_\xi - \text{a.s.}$  and ?? shows that  $\mathbb{E}_\xi[\mathbb{1}_A | \mathcal{F}_n] \xrightarrow{\mathbb{P}_\xi - \text{a.s.}} \mathbb{1}_A \quad \mathbb{P}_\xi - \text{a.s.}$  which implies that  $c \in \{0, 1\}$ .

(ii)  $\Rightarrow$  (i): Let  $h$  be a bounded harmonic function and  $\xi \in \mathbb{M}_1(\mathcal{X})$ . The random variable  $Y = \limsup_{n \rightarrow \infty} h(X_n)$  is invariant and under (ii), there exists a constant  $c < \infty$  (possibly depending on  $\xi$ ) such that  $\limsup_{n \rightarrow \infty} h(X_n) = c \quad \mathbb{P}_\xi - \text{a.s.}$  By Proposition 5.17  $\mathbb{E}_x[Y] = c = h(x)$  for all  $x \in \mathbf{X}$ . Therefore,  $h$  is constant which shows (i). ■

If we wish to obtain the Law of Large Numbers for a particular Markov chain by applying Theorem 5.14, we have to check the ergodicity assumption. It is therefore convenient to have sufficient conditions ensuring ergodicity.

We now give a sufficient condition for  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  to be ergodic expressed in terms of absorbing sets.

**Lemma 5.20** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  admitting an invariant probability measure  $\pi$ . If for all absorbing sets  $B \in \mathcal{X}$ ,  $\pi(B) \in \{0, 1\}$ , then the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic.*

PROOF. Let  $A \in \mathcal{S}$  and define  $h(x) = \mathbb{E}_x[\mathbb{1}_A]$  and  $B = \{x \in X : h(x) = 1\}$ . By Proposition 5.17-(i),  $h$  is a nonnegative harmonic function bounded by 1. For any  $x \in B$  we have  $\mathbb{E}_x[h(X_1)] = Ph(x) = h(x) = 1$ , which implies  $\mathbb{P}_x(h(X_1) = 1) = 1$ . Therefore for any  $x \in B$ , we get  $\mathbb{P}_x(X_1 \in B) = \mathbb{P}_x(h(X_1) = 1) = 1$ . Therefore  $B$  is absorbing and hence, under the stated assumption, we have  $\pi(B) \in \{0, 1\}$ .

By Proposition 5.17 (iii), we know that  $\mathbb{P}_\pi(\mathbb{E}_{X_0}[\mathbb{1}_A] = \mathbb{1}_A) = 1$  which implies that  $\mathbb{P}_\pi(h(X_0) \in \{0, 1\}) = 1$ . This yields

$$\begin{aligned} \mathbb{P}_\pi(A) &= \mathbb{E}_\pi[\mathbb{E}_{X_0}[\mathbb{1}_A]] = \int_X \pi(dx) h(x) \\ &= \int_X \pi(dx) \mathbb{1}_{\{h(x) = 1\}} = \int_X \pi(dx) \mathbb{1}_B(x) = \pi(B). \end{aligned}$$

Thus  $\mathbb{P}_\pi(A) \in \{0, 1\}$  and  $\mathcal{S}$  is trivial for  $\mathbb{P}_\pi$ . ■

It turns out that the sufficient condition in Lemma 5.20 is also a necessary condition. Before showing the necessary part, we first draw an easy and useful consequence of Lemma 5.20.

**Theorem 5.21.** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  admitting a unique invariant probability measure  $\pi$ . The dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic.*

PROOF. By Proposition 5.16, since  $P$  has unique invariant probability  $\pi$ , for every absorbing set  $B$ ,  $\pi(B) \in \{0, 1\}$ . We conclude by Lemma 5.20. ■

The uniqueness of the invariant probability measure is a sufficient but not a necessary condition for ergodicity as illustrated in Exercise 5.35. Comparing with Lemma 5.20, the following Lemma goes one step further. When the dynamical system is not ergodic, the state space  $X$  contains not only one but at least two disjoint absorbing sets which are not trivial with respect to  $\pi$ .

**Lemma 5.22** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  admitting an invariant probability measure  $\pi$ . If the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is not ergodic, then, there exist two disjoint absorbing sets  $B$  and  $B'$  in  $\mathcal{X}$  such that  $\pi(B) = 1 - \pi(B') \in (0, 1)$  and  $\pi_B(\cdot) = \pi(B \cap \cdot)/\pi(B)$  and  $\pi_{B'}(\cdot) = \pi(B' \cap \cdot)/\pi(B')$  are invariant probability measures.*

PROOF. Since the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is not ergodic, there exists  $A \in \mathcal{S}$  such that  $\mathbb{P}_\pi(A) = \alpha \in (0, 1)$ . As  $\mathcal{S}$  is a  $\sigma$ -field, we also have  $A^c \in \mathcal{S}$ . Define  $B = \{x \in X : \mathbb{E}_x[\mathbb{1}_A] = 1\}$  and  $B' = \{x \in X : \mathbb{E}_x[\mathbb{1}_{A^c}] = 1\}$ . As noted in the proof of Lemma 5.20 (see also Exercise 5.31), the sets  $B$  and  $B'$  are absorbing and

$$\pi(B) = \mathbb{P}_\pi(A) = 1 - \mathbb{P}_\pi(A^c) = 1 - \pi(B') \in (0, 1).$$

By Proposition 5.16,  $\pi_B$  and  $\pi_{B'}$  are invariant probability measures. ■

Without ergodicity assumption, the generalized version of Birkhoff's ergodic theorem in Theorem 5.8 shows that the normalized partial sums still converges but the limit is a random variable which is not necessarily almost surely constant (see also an illustration in Exercise 5.35). In the context of Markov chains this limit turns out to be a function of  $X_0$ . More precisely we have the following theorem.

**Proposition 5.23** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  admitting an invariant probability measure  $\pi$  and let  $Y \in L^1(\mathbb{P}_\pi)$ . Let  $Z$  be a version of the conditional expectation  $\mathbb{E}_\pi[Y | \mathcal{S}]$ , i.e.  $Z = \mathbb{E}_\pi[Y | \mathcal{S}]$   $\mathbb{P}_\pi$ -a.s. Then, there exists a set  $S \in \mathcal{X}$ , such that  $\pi(S) = 1$  and for each  $x \in S$ ,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_x[Z] \quad \mathbb{P}_x - \text{a.s.} \quad (5.10)$$

PROOF. Define  $\phi(x) = \mathbb{E}_x[Z]$ . It follows from Proposition 5.17 (iii) that  $\mathbb{E}_\pi[Y | \mathcal{J}] = \phi(X_0)$ ,  $\mathbb{P}_\pi - \text{a.s.}$  Hence, Theorem 5.8 yields

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_\pi[Y | \mathcal{J}] = \phi(X_0) \quad \mathbb{P}_\pi - \text{a.s.}$$

Set  $A = \{\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} Y \circ \theta_k = \phi(X_0)\}$ . The previous relation implies  $\mathbb{P}_\pi(A) = 1$ , i.e.  $\int \pi(dx) \mathbb{P}_x(A) = 1$ . Since  $\mathbb{P}_x(A) \leq 1$  for all  $x \in X$ , this implies that  $\mathbb{P}_x(A) = 1$  for  $\pi$ -almost all  $x \in X$ . Setting  $S = \{x \in X : \mathbb{P}_x(A) = 1\}$  concludes the proof. ■

In Proposition 5.23, the limit  $\mathbb{E}_x[Z]$  in (5.10) is expressed in terms of  $Z$ , a version of the conditional expectation  $\mathbb{E}_\pi[Y | \mathcal{J}]$ . If we choose another version of  $\mathbb{E}_\pi[Y | \mathcal{J}]$ , say  $Z'$ , under  $\mathbb{P}_\pi$ , then obviously,  $Z = Z'$   $\mathbb{P}_\pi - \text{a.s.}$  but we do not necessarily have  $\mathbb{E}_x[Z'] = \mathbb{E}_x[Z]$   $\mathbb{P}_x - \text{a.s.}$  since without additional assumption,  $\mathbb{P}_x$  is not necessarily dominated by  $\mathbb{P}_\pi$ . The situation is different when the dynamical system is ergodic since the limit is then  $\mathbb{P}_\pi - \text{a.s.}$  constant.

**Theorem 5.24 (Birkhoff's Theorem for Markov chains).** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  and assume that  $P$  admits an invariant probability measure  $\pi$  such that  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic. Let  $Y \in L^1(\mathbb{P}_\pi)$ . Then, for  $\pi$ -almost all  $x \in X$ ,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_\pi[Y] \quad \mathbb{P}_x - \text{a.s.}$$

PROOF. Since  $\pi$  is ergodic, the invariant  $\sigma$ -field  $\mathcal{J}$  is trivial for  $\mathbb{P}_\pi$ . This implies  $\mathbb{E}_\pi[Y | \mathcal{J}] = \mathbb{E}_\pi[Y]$   $\mathbb{P}_\pi - \text{a.s.}$  ■

**Theorem 5.25.** *Let  $P$  a Markov kernel on  $X \times \mathcal{X}$ . If  $\pi_1$  and  $\pi_2$  are distinct invariant probability measures such that  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\pi_1}, \theta)$  and  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\pi_2}, \theta)$  are ergodic, then  $\pi_1$  and  $\pi_2$  are mutually singular.*

PROOF. Note first that, if  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic and  $f \in \mathbb{F}_b(X)$ , then applying Theorem 5.24 to the random variable  $Y = f(X_0)$  and the dominated convergence theorem, we obtain that there exists a set  $S \in \mathcal{X}$ , such that  $\pi(S) = 1$  and for all  $x \in S$ ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} P^k f(x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} \mathbb{E}_x[f(X_k)] = \pi(f).$$

Now assume that  $\pi_1$  and  $\pi_2$  are different invariant probability measures such that  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\pi_1}, \theta)$  and  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\pi_2}, \theta)$  are ergodic. Let  $C \in \mathcal{X}$  such that  $\pi_1(C) \neq \pi_2(C)$  and set, for  $i = 1, 2$ ,

$$S_i = \left\{ x \in X : \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} P^k \mathbb{1}_C(x) = \pi_i(C) \right\},$$

We have  $S_1 \cap S_2 = \emptyset$ ,  $\pi_1(S_1) = 1$  and  $\pi_2(S_2) = 1$ , which means that  $\pi_1$  and  $\pi_2$  are mutually singular. ■

We have now all the tools for getting a necessary and sufficient condition for the dynamical system to be ergodic.

**Theorem 5.26.** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  admitting an invariant probability measure  $\pi$ . The dynamical system  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic if and only if for all absorbing sets  $B \in \mathcal{X}$ ,  $\pi(B) \in \{0, 1\}$ .*

PROOF. The sufficient condition follows from Lemma 5.20. We now consider an ergodic dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  and we let  $B$  be an absorbing set. Assume first that  $\pi(B) > 0$ . Then, by Proposition 5.16,  $\tilde{\pi}_B(\cdot) = \pi(B \cap \cdot) / \pi(B)$  is an invariant probability measure. Moreover, note that for all  $A \in \mathcal{X}^{\otimes \mathbb{N}}$ ,

$$\mathbb{P}_{\tilde{\pi}_B}(A) = \int \pi(dx) \mathbb{P}_x(A) \frac{\mathbb{1}_B(x)}{\pi(B)} \leq \frac{\mathbb{P}_\pi(A)}{\pi(B)}.$$

Combining with the ergodicity of the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  we deduce that any invariant set  $A \in \mathcal{I}$  satisfies either  $0 = \mathbb{P}_\pi(A) = \mathbb{P}_{\tilde{\pi}_B}(A)$  or  $0 = \mathbb{P}_\pi(A^c) = \mathbb{P}_{\tilde{\pi}_B}(A^c)$ . The dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\tilde{\pi}_B}, \theta)$  is therefore ergodic and by Theorem 5.25,  $\tilde{\pi}_B = \pi$  since they are not mutually singular. This implies that  $\pi(B^c) = \tilde{\pi}_B(B^c) = 0$ . Finally,  $\pi(B) \in \{0, 1\}$ , which concludes the proof. ■

In Proposition 5.23 and Theorem 5.24, the law of large numbers is obtained under  $\mathbb{P}_x$  for all  $x$  belonging to a set  $S$  such that  $\pi(S) = 1$  and which may depend on the random variable  $Y$  under consideration. This is unsatisfactory since it does not tell if the LLN holds for a given  $x \in X$  or more generally for a given initial distribution  $\xi$ . We now give a criterion to obtain the LLN when the chain does not start from stationarity. Recall that the total variation distance between two probability measures  $\mu, \nu \in \mathbb{M}_1(\mathcal{X})$  is defined by

$$\|\mu - \nu\|_{TV} = \sup_{h \in \mathbb{F}_b(X), |h|_\infty \leq 1} |\mu(h) - \nu(h)|.$$

More details and basic properties on the total variation distance are given in ??.

**Proposition 5.27** *Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  admitting an invariant probability measure  $\pi$ . If the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic and if  $\xi \in \mathbb{M}_1(\mathcal{X})$  is such that  $\lim_{n \rightarrow \infty} \|n^{-1} \sum_{k=1}^n \xi P^k - \pi\|_{TV} = 0$  then for all  $Y \in \mathcal{L}^1(\mathbb{P}_\pi)$ ,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_\pi[Y] \quad \mathbb{P}_\xi - \text{a.s.}$$

PROOF. Set  $A = \{\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_\pi[Y]\}$ . Since  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic, we already know that  $\mathbb{P}_\pi(A) = 1$ . We show that  $\mathbb{P}_\xi(A) = 1$ . Define the function  $h$  by  $h(x) = \mathbb{E}_x[\mathbb{1}_A]$ . Since  $A \in \mathcal{I}$ , Proposition 5.17 implies that  $h$  is harmonic. Then, for all  $n \in \mathbb{N}$ ,  $n^{-1} \sum_{k=1}^n \xi P^k h = \xi(h)$ . Moreover, noting that  $h \leq 1$ , we have by assumption,  $\lim_{n \rightarrow \infty} n^{-1} \sum_{k=1}^n \xi P^k(h) = \pi(h)$ . Thus,  $\xi(h) = \pi(h)$  and

$$\mathbb{P}_\xi(A) = \int_X \xi(dx) \mathbb{E}_x[\mathbb{1}_A] = \xi(h) = \pi(h) = \mathbb{P}_\pi(A) = 1.$$

The condition  $\lim_n \|\xi P^n - \pi\|_{TV} = 0$  is not mandatory for having a Law of Large Numbers under  $\mathbb{P}_\xi$ . In some situations, one can get the same result without any straightforward information in the decrease of  $\|\xi P^n - \pi\|_{TV}$  toward 0. This is the case for example for Metropolis-Hastings kernels as illustrated in Exercise 5.36. Another illustration can be found in Exercise 5.41 where the Law of Large Numbers is extended to different initial distributions in the case where  $(X, d)$  is a complete separable metric space.

## 5.3 Exercises

**Exercise 5.28.** Let  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be a dynamical system. Show that  $\mathcal{I} \neq \bigcap_{k \geq 0} \sigma(X_l, l > k)$ .

**Exercise 5.29.** Let  $(\Omega, \mathcal{B}, \mathbb{P})$  be a probability space and  $\theta : \Omega \rightarrow \Omega$  be a measurable transformation. Let  $\mathcal{B}_0$  a family of sets, stable under finite intersection and generating  $\mathcal{B}$ . If for all  $B \in \mathcal{B}_0$ ,  $\mathbb{P}[\theta^{-1}(B)] = \mathbb{P}(B)$ , then  $T$  is measure-preserving.

**Exercise 5.30.** Let  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be a dynamical system. Let  $Y$  be a  $\mathbb{R}$ -valued random variable such that  $\mathbb{E}[Y^+] < \infty$ . Show that for all  $k \geq 0$ ,

$$\mathbb{E} [Y \circ T^k | \mathcal{J}] = \mathbb{E} [Y | \mathcal{J}] \quad \mathbb{P} - \text{a.s.}$$

**Exercise 5.31.** Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$  and let  $(\Omega, \mathcal{B})$  be the canonical space. For  $A \in \mathcal{J}$ , define  $B = \{x \in X : \mathbb{P}_x(A) = 1\}$ .

1. Show that  $B$  is absorbing.
2. Let  $\pi$  be an invariant probability. Show that  $\pi(A) = \mathbb{P}_\pi(B)$ .

**Exercise 5.32.** The following exercise provides the converse of Theorem 5.14. Let  $P$  be a Markov kernel on  $X \times \mathcal{X}$ . Let  $\pi$  be a probability measure,  $\pi \in \mathbb{M}_1(\mathcal{X})$ . Assume that for all  $f \in \mathbb{F}_b(X)$ , we get

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} f(X_k) = \pi(f) \quad \mathbb{P}_\pi - \text{a.s.}$$

1. Show that  $\pi$  is invariant.
2. Let  $A \in \mathcal{J}$  and set  $B = \{x \in X : \mathbb{P}_x(A) = 1\}$ . Show that

$$\mathbb{1}_A = \mathbb{P}_{X_0}(A) = \mathbb{1}_B(X_0) \quad \mathbb{P}_\pi - \text{a.s.}$$

and that for all  $k \in \mathbb{N}$ ,  $\mathbb{1}_A = \mathbb{1}_B(X_0) = \dots = \mathbb{1}_B(X_k)$   $\mathbb{P}_\pi - \text{a.s.}$ .

3. Show that the dynamical system  $(X^\mathbb{N}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic.

**Exercise 5.33.** In this exercise, we prove various extensions of Birkhoff's ergodic theorem. Let  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be a dynamical system. In the first two questions, we assume that the dynamical system is ergodic.

1. Let  $Y$  be nonnegative random variable such that  $\mathbb{E}[Y] = \infty$ . Show that

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \infty \quad \mathbb{P} - \text{a.s.}$$

[Hint: use Corollary 5.11 with  $Y_M = Y \wedge M$  and let  $M$  tends to infinity.]

2. Let  $Y$  be a random variable such that  $\mathbb{E}[Y^+] < \infty$ . Show that

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \mathbb{E}[Y] \quad \mathbb{P} - \text{a.s.}$$

3. In what follows, we do not assume any ergodicity of the dynamical system  $(\Omega, \mathcal{B}, \mathbb{P}, T)$ . Let  $Y$  be a nonnegative random variable. Set

$$A = \left\{ \lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \mathbb{E}[Y | \mathcal{J}] \right\}$$

Let  $M > 0$ . Using Theorem 5.8 with  $Y \mathbb{1}_{\{\mathbb{E}[Y | \mathcal{J}] \leq M\}}$ , show that on  $\{\mathbb{E}[Y | \mathcal{J}] \leq M\}$

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \mathbb{E}[Y | \mathcal{J}] \quad \mathbb{P} - \text{a.s.}$$

Deduce that  $\mathbb{P}(A^c \cap \{\mathbb{E}[Y | \mathcal{J}] < \infty\}) = 0$ . Moreover, show that for all  $M > 0$ ,

$$\liminf_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k \geq \mathbb{E}[Y \wedge M | \mathcal{J}] \quad \mathbb{P} - \text{a.s.}$$

Deduce that  $\mathbb{P}(A^c \cap \{\mathbb{E}[Y | \mathcal{J}] = \infty\}) = 0$ .

4. Let  $Y$  be a random variable such that  $\mathbb{E}[Y^+] < \infty$ , show that

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ T^k = \mathbb{E}[Y | \mathcal{I}] \quad \mathbb{P} - \text{a.s.}$$

[Hint: write  $Y = Y^+ - Y^-$  and use the previous question with  $Y^-$ ]

**Exercise 5.34.** Let  $P$  be the kernel on  $X \times \mathcal{X}$  defined by for all  $(x, A) \in X \times \mathcal{X}$ ,  $P(x, A) = \delta_x(A)$ . Find all the probability measures  $\pi \in \mathbb{M}_1(\mathcal{X})$  such that  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_\pi, \theta)$  is ergodic.

**Exercise 5.35.** Let  $\mu_0$  (resp.  $\mu_1$ ) be a probability measure on  $\mathbb{R}^+$  (resp.  $\mathbb{R}^- \setminus \{0\}$ ). Let  $P$  be a Markov kernel on  $\mathbb{R} \times \mathcal{B}(\mathbb{R})$  defined by  $P(x, \cdot) = \mu_0$  if  $x \geq 0$  and  $P(x, \cdot) = \mu_1$  otherwise. Set for all  $\alpha \in (0, 1)$ ,  $\mu_\alpha = (1 - \alpha)\mu_0 + \alpha\mu_1$ .

1. Show that for all  $\alpha \in (0, 1)$ ,  $\mu_\alpha$  is an invariant probability measure but the dynamical system  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\mu_\alpha}, \theta)$  is not ergodic.
2. Show that for  $\alpha \in \{0, 1\}$ ,  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\mu_\alpha}, \theta)$  is ergodic.
3. Let  $f \in L^1(\mu_0) \cap L^1(\mu_1)$ . Find a function  $\phi$  such that for all  $x \in \mathbb{R}$ ,

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} f(X_k) = \phi(X_0) \quad \mathbb{P}_{\mu_\alpha} - \text{a.s.}$$

For all  $\alpha \in (0, 1)$ , find a version of  $\mathbb{E}_{\mu_\alpha}[f | \mathcal{I}]$ .

**Exercise 5.36.** In this exercise, we will find a sufficient condition for a Metropolis-Hastings kernel to satisfy the law of large numbers starting from any initial distribution. We make use of the notation of Section 4.1.1. Let  $\pi$  be target distribution on a measurable space  $(X, \mathcal{X})$  and assume that  $\pi$  has a positive density  $h$  with respect to a measure  $\mu \in \mathbb{M}_+(\mathcal{X})$ . Let  $Q$  be a proposal kernel on  $X \times \mathcal{X}$  and assume that  $Q$  has a positive kernel density  $y \mapsto q(x, y)$  with respect to  $\mu$ . The Metropolis-Hasting kernel  $P$  is then defined by (4.4).

1. Show that

$$P(x, A) \geq \int_A \alpha(x, y) q(x, y) \mu(dy).$$

Deduce that  $\pi$  is the unique invariant probability of  $P$ .

2. Let  $A \in \mathcal{I}$  and assume that  $\mathbb{P}_\pi(A) = 0$ . Set  $\phi(x) = \mathbb{P}_x(A)$ . For all  $x \in X$ , show that

$$\phi(x) = P\phi(x) = \int \frac{\alpha(x, y) q(x, y)}{h(y)} \pi(dy) \phi(y) + \bar{\alpha}(x) \phi(x).$$

and deduce that  $\mathbb{P}_x(A) = 0$ .

3. Let  $\xi \in \mathbb{M}_1(\mathcal{X})$ . Deduce from the previous question that for all random variables  $Y \in L^1(\mathbb{P}_\pi)$ ,

$$\lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} Y \circ \theta_k = \mathbb{E}_\pi[Y] \quad \mathbb{P}_\xi - \text{a.s.}$$

**Exercise 5.37.** Let  $P$  a Markov kernel on  $X \times \mathcal{X}$  admitting an invariant probability measure  $\pi_1$  such that  $(X^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}}, \mathbb{P}_{\pi_1}, \theta)$  is ergodic. Let  $\mu$  be another invariant probability measure.

1. Show that  $\pi_1 \wedge \mu$  is an invariant finite measure.
2. If  $\pi_1 \wedge \mu(X) \neq 0$ , show, using Theorem 5.25, that  $\pi_1 \wedge \mu / \pi_1 \wedge \mu(X) = \pi_1$ .
3. Show that there exists an invariant probability measure  $\pi_2$  satisfying:  $\pi_1$  and  $\pi_2$  are mutually singular and there exists  $\alpha \in [0, 1]$  such that  $\mu = \alpha\pi_1 + (1 - \alpha)\pi_2$ .

**Exercise 5.38.** Let  $\{X_n, n \in \mathbb{Z}\}$  be a canonical stationary Markov chain on  $(X^{\mathbb{Z}}, \mathcal{X}^{\otimes \mathbb{Z}})$ . We set

$$\overline{\mathcal{F}}_{-\infty}^0 = \overline{\sigma(X_k, k \leq 0)}^P, \quad \overline{\mathcal{F}}_0^\infty = \overline{\sigma(X_k, k \geq 0)}^P.$$

We consider an invariant bounded random variable  $Y$ .

- (i) Show that  $Y$  is  $\overline{\mathcal{F}_0^\infty}$  and  $\overline{\mathcal{F}_{-\infty}^0}$  measurable.
- (ii) Deduce from the previous question that  $Y = \mathbb{E}[Y | X_0] \quad \mathbb{P} - \text{a.s.}$
- (iii) Show that the previous identity holds true for all  $\mathbb{P}$ -integrable or positive invariant random variable  $Y$ .

The following exercises deal with subadditive sequences. A sequence of random variables  $\{Y_n, n \in \mathbb{N}^*\}$  is said to be subadditive for the dynamical system  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  if for all  $(n, p) \in \mathbb{N}^*$ ,  $Y_{n+p} \leq Y_n + Y_p \circ T^n$ . The sequence is said to be additive if for all  $(n, p) \in \mathbb{N}^*$ ,  $Y_{n+p} = Y_n + Y_p \circ T^n$ .

**Exercise 5.39 (Fekete Lemma).** Consider  $\{a_n, n \in \mathbb{N}^*\}$ , a sequence in  $[-\infty, \infty)$  such that, for all  $(m, n) \in \mathbb{N}^* \times \mathbb{N}^*$ ,  $a_{n+m} \leq a_n + a_m$ . Then,

$$\lim_{n \rightarrow \infty} \frac{a_n}{n} = \inf_{m \in \mathbb{N}^*} \frac{a_m}{m} ;$$

in other words, the sequence  $\{n^{-1}a_n, n \in \mathbb{N}^*\}$  either converges to its lower bounds or diverges to  $-\infty$ .

**Exercise 5.40.** Let  $(\Omega, \mathcal{B}, \mathbb{P}, T)$  be a dynamical system and  $\{Y_n, n \in \mathbb{N}^*\}$  be a subadditive sequence of functions such that  $\mathbb{E}[Y_1^+] < \infty$ . Show that for any  $n \in \mathbb{N}^*$ ,  $\mathbb{E}[Y_n^+] \leq n\mathbb{E}[Y_1^+] < \infty$  and

$$\lim_{n \rightarrow \infty} n^{-1}\mathbb{E}[Y_n] = \inf_{n \in \mathbb{N}^*} n^{-1}\mathbb{E}[Y_n] , \quad (5.11)$$

$$\lim_{n \rightarrow \infty} n^{-1}\mathbb{E}[Y_n | \mathcal{I}] = \inf_{n \in \mathbb{N}^*} n^{-1}\mathbb{E}[Y_n | \mathcal{I}] , \quad \mathbb{P} - \text{a.s.} \quad (5.12)$$

where  $\mathcal{I}$  is the invariant  $\sigma$ -field [Hint: use Exercise 5.30].

**Exercise 5.41.** Let  $P$  be a Markov kernel on a complete separable metric space  $(X, d)$  which admits a unique invariant probability measure  $\pi$ . We assume that there exists a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and a stochastic process  $\{(X_n, X'_n), n \in \mathbb{N}\}$  such that  $\{X_n, n \in \mathbb{N}\}$  and  $\{X'_n, n \in \mathbb{N}\}$  are Markov chains with kernel  $P$  and initial distribution  $\xi \in \mathbb{M}_1(\mathcal{X})$  and  $\pi$ , respectively. Assume that  $d(X_n, X'_n) \xrightarrow{\mathbb{P}-\text{a.s.}} 0$ .

We recall the Parthasaraty's theorem (see (2, Theorem 6.6)): Then there exists a countable set  $H$  of bounded continuous functions such that for all  $\{\mu, \mu_n, n \geq 1\} \subset \mathbb{M}_1(X)$ , the following assertions are equivalent:

- (i)  $\mu_n$  converges weakly to  $\mu$ .
- (ii) For all  $h \in H$ ,  $\lim_{n \rightarrow \infty} \mu_n(h) = \mu(h)$ .

For this set  $H$ , define the event

$$A = \left\{ \omega \in \Omega : \forall h \in H, \lim_{n \rightarrow \infty} n^{-1} \sum_{k=0}^{n-1} h(X'_k(\omega)) = \pi(h) \right\} .$$

1. Show that  $\mathbb{P}(A) = 1$ .
2. Deduce that there exists a set  $\tilde{\Omega}$  such that  $\mathbb{P}(\tilde{\Omega}) = 1$  and for all bounded continuous functions  $h$  on  $X$  and all  $\tilde{\omega} \in \tilde{\Omega}$ ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} h(X_k(\tilde{\omega})) = \pi(h) .$$

Let  $V$  be a nonnegative and uniformly continuous such that  $\pi(V) < \infty$ .

3. Show that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} V(X_k) = \pi(V) , \quad \mathbb{P} - \text{a.s.}$$

4. Show that there exists  $\tilde{\Omega}$  such that for all  $\omega \in \tilde{\Omega}$  and all continuous functions  $f$  such that  $\sup_{x \in X} |f(x)|/V(x) < \infty$ ,



$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} f(X_k(\omega)) = \pi(f), \quad \mathbb{P} - \text{a.s.}$$

## 5.4 Bibliographical notes

Ergodic theory is a very important area of probability theory which has given rise to a great deal of work. The interested reader will find an introduction to this field in the books ? and ?. The application of ergodic theory to Markov chains is a very classic subject. A detailed study of the ergodic theory of Markov chains can be found in (?, Chapter 4) and (?, Chapter 2); These books contain many references to works on this subject that began in the early 1960s.

The proof of the Birkhoff Theorem (Theorem 5.8) is borrowed from unpublished notes by B. Delyon and extended to possibly non-ergodic dynamical systems. This approach is closely related to the very short proof of the Law of Large Numbers for i.i.d. random variables written by Jacques Neveu in his unpublished lecture notes at Ecole Polytechnique. Several other proofs of the ergodic theorem are also given in ?.

Theorem 5.21 is essentially borrowed from (?, Proposition 2.4.3) even if the statements of the two results are slightly different.



# Chapter 6

## Pseudo Marginal Monte Carlo methods and applications



# Chapter 7

## Hamiltonian Monte Carlo methods

### Contents

7.1	MH with deterministic moves .....	97
7.2	Hamiltonian dynamics .....	98
7.2.1	Level sets .....	98
7.2.2	Involution .....	98
7.3	The leapfrog integrator .....	99

Let  $\pi(q) \propto e^{-U(q)}$  be a target distribution. We may consider  $\pi$  as the marginal of the extended target

$$\Pi(q, p) \propto \exp \{ -U(q) - p^2/2 \} .$$

We make use of the following terminology

- $q \in \mathbb{R}^d$  is the position and  $U(q)$  is the potential energy.
- $p \in \mathbb{R}^d$  is the momentum and  $K(p) = p^T p/2$  is the kinetic energy.
- $H(q, p) = U(q) + K(p)$  is called the Hamiltonian.

### 7.1 MH with deterministic moves

If we want a deterministic move to target  $\Pi$  in a Metropolis Hasting algorithm, then the proposal kernel writes:  $Q(x, dy) = \delta_{\varphi(x)}(dy)$  and the acceptance probability is

$$\frac{d\mu}{dv}(x, y) \wedge 1 \quad \text{where} \quad \mu(dx dy) = \Pi(dy) \delta_{\varphi(y)}(dx), \quad v(dx dy) = \Pi(dx) \delta_{\varphi(x)}(dy)$$

but this is only possible if  $\mu$  and  $v$  are equivalent (ie dominated one by the other and conversely). To get that, write

$$\begin{aligned} \int_A \mu(dx dy) &= \int \Pi(du) \mathbb{1}_A(\underbrace{\varphi(u)}_v, \underbrace{u}_{\varphi^{-1}(v)}) \\ &= \int \Pi \circ \varphi^{-1}(dv) \mathbb{1}_A(v, \varphi^{-1}(v)) = \int \frac{d\Pi \circ \varphi^{-1}}{d\Pi}(v) \Pi(dv) \mathbb{1}_A(v, \varphi^{-1}(v)) = \int_A \frac{d\Pi \circ \varphi^{-1}}{d\Pi}(x) v(dx dy) \end{aligned}$$

only if  $\varphi^{-1}(v) = \varphi(v)$  that is  $\varphi$  is an **involution**.

**Remark 7.1** (i) For any involution, you can get a Metropolis Hasting with a theoretical expression of the acceptance probability as

$$\frac{d\Pi \circ \varphi^{-1}}{d\Pi}(x) \wedge 1$$

but the ideal HMC goes one step further since we can show that this is equal to 1. To get this, if we work on  $\mathbb{R}^d$  and if  $\Pi$  has density wrt the Lebesgue measure that we still denote  $\Pi$ , we get

$$\frac{d\Pi \circ \varphi^{-1}}{d\Pi}(x) = \frac{\Pi(\varphi^{-1}(x))}{\Pi(x)} \left| \frac{\partial \varphi^{-1}(x)}{\partial x} \right|$$

where the second term  $\left| \frac{\partial \varphi^{-1}(x)}{\partial x} \right|$  is the Jacobian of the transformation  $\varphi$ . To get 1 in the acceptance probability, we can impose that the two terms are equal to 1. The first term is one if the involution stays on the same level set and the second is one if the involution is volume invariant. If for example the involution only keeps the volume then the Radon Nikodym simplifies to

$$\frac{\Pi(\varphi^{-1}(x))}{\Pi(x)} = \frac{\Pi(\varphi(x))}{\Pi(x)} \quad \text{since} \quad \varphi \circ \varphi = \text{I}$$

- (ii) The second point is that if we only use the involution, then after two steps we land up to the initial state... Not very interesting... Therefore, this deterministic transition is often combined with another move that is not deterministic.

## 7.2 Hamiltonian dynamics

### 7.2.1 Level sets

If we now let  $(p, q)$  depend on a real parameter  $t$  and we impose to stay on a level set of  $H$ , we get:

$$\frac{dH(q_t, p_t)}{dt} = 0 = \sum_{i=1}^d \frac{\partial H(q_t, p_t)}{\partial q_t^i} \frac{dq_t^i}{dt} + \frac{\partial H(q_t, p_t)}{\partial p_t^i} \frac{dp_t^i}{dt}.$$

This gives the idea of using the following dynamics: for all  $i \in \{1, \dots, d\}$ ,

$$\begin{aligned} \frac{\partial H}{\partial q_t^i}(q_t, p_t) &= \frac{\partial U(q_t)}{\partial q_t^i} = -\frac{dp_t^i}{dt} \\ \frac{\partial H}{\partial p_t^i}(q_t, p_t) &= \frac{\partial K(p_t)}{\partial p_t^i} = p_t^i = \frac{dq_t^i}{dt} \end{aligned}$$

The very last equation leads to the interpretation of  $p_t^i$  as a speed since it is the derivative of the position wrt time. Note  $\phi_t(q, p) = (q_t, p_t)$  the (deterministic) position and momentum at time  $t$  when  $(q_s, p_s)$  follows the Hamiltonian dynamics.

### 7.2.2 Involution

Denote  $s(q, p) = (q, -p)$  and set  $f_T = s \circ \phi_T$ . We now show that  $f_T$  is an involution. Indeed, write  $f_T(q, p) = (q_T, -p_T)$ . To see what we obtain by applying again  $f_T$ , set  $\tilde{q}_t = q_{T-t}$  and  $\tilde{p}_t = -p_{T-t}$  so that  $(\tilde{q}_0, \tilde{p}_0) = f_T(q, p)$ . Then,

$$\begin{aligned}\frac{d\tilde{q}_t^i}{dt} &= \frac{dq_{T-t}^i}{dt} = - \left. \frac{dq_s^i}{ds} \right|_{s=T-t} = -p_{T-t}^i = \tilde{p}_t^i \\ \frac{d\tilde{p}_t^i}{dt} &= - \frac{dp_{T-t}^i}{dt} = - \left. \frac{dp_s^i}{ds} \right|_{s=T-t} = - \frac{\partial U(q_{T-t})}{\partial q_{T-t}^i} = - \frac{\partial U(\tilde{q}_{T-t})}{\partial \tilde{q}_{T-t}^i}\end{aligned}$$

Finally, the process  $(\tilde{q}_t, \tilde{p}_t)$  follows the Hamiltonian dynamics so that  $f_T(\tilde{q}_0, \tilde{p}_0) = (\tilde{q}_T, -\tilde{p}_T)$  and by definition this quantity is equal to  $(q_0, p_0)$ . We finally obtain that  $f_T$  is an involution.

### 7.3 The leapfrog integrator

$$\begin{aligned}p_{k+1/2} &= p_k - (h/2)\nabla U(q_k) \\ q_{k+1} &= q_k - hp_{k+1/2} \\ p_{k+1} &= p_{k+1/2} - (h/2)\nabla U(q_{k+1})\end{aligned}$$





# Chapter 8

## Exercises

**Exercise 8.1.** Let  $\{Z_t, t \in \mathbb{N}^*\}$  be an i.i.d. Bernoulli sequence such that  $\mathbb{P}(Z_t = 1) = p = 1 - \mathbb{P}(Z_t = 0)$ , and  $p \in (0, 1)$ . Define  $\{X_t, t \in \mathbb{N}\}$  by  $X_0 = 0$  and  $X_t = X_{t-1} + Z_t = \sum_{s=1}^t Z_s$ , for  $t \geq 1$ . Denote by  $\mathcal{F}_t^X = \sigma(X_s, 0 \leq s \leq t)$  the  $\sigma$ -algebra generated by the random variables  $\{X_0, \dots, X_t\}$ .

(a) Show that

$$\mathbb{P}(X_{t+1} = x_{t+1} \mid X_t = x_t, \dots, X_1 = x_1) = \mathbb{P}(X_{t+1} = x_{t+1} \mid X_t = x_t) = \begin{cases} p & \text{if } x_{t+1} - x_t = 1, \\ 1 - p & \text{if } x_{t+1} - x_t = 0, \\ 0 & \text{otherwise.} \end{cases}$$

(b) Consider  $A = \{X_1 = 1, X_2 = 1\}$  and  $B = \{X_2 = 1, X_3 = 1\}$ . Show that  $A = \{Z_1 = 1, Z_2 = 0\}$  and  $B = \{(Z_1 = 0, Z_2 = 1, Z_3 = 0) \cup (Z_1 = 1, Z_2 = Z_3 = 0)\}$ . Show that  $A$  and  $B$  are not independent.

(c) Show that  $A$  and  $B$  are conditionally independent given  $X_1$ .  $\mathbb{P}(A \cap B \mid X_2 = 1) = (1 - p)/2 = \mathbb{P}(A \mid X_2 = 1)\mathbb{P}(B \mid X_2 = 1)$ .

**Exercise 8.2.** Let  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  be a Markov chain.

(a) Show that for all  $0 \leq r \leq t$  and  $f \in \mathbb{F}_b(X)$

$$\mathbb{E}[f(X_{t+1}) \mid \mathcal{F}_t] = \mathbb{E}[f(X_{t+1}) \mid \sigma(X_r, \dots, X_t)] = \mathbb{E}[f(X_{t+1}) \mid \mathcal{F}_t],$$

(b) Consider the property  $(\mathcal{P}_n)$ : for all  $Y = \prod_{k=0}^n g_k(X_{t+k})$  where  $g_k \in \mathbb{F}_b(X)$ ,  $\mathbb{E}[Y \mid \mathcal{F}_t] = \mathbb{E}[Y \mid X_t]$ ,  $\mathbb{P}$ -a.s. Show that  $\mathcal{P}_0$  is satisfied.

(c) Assume that  $(\mathcal{P}_n)$  is satisfied. Then, for any  $g_k \in \mathbb{F}_b(X)$ , prove that

$$\mathbb{E}[g_0(X_t) \dots g_n(X_{t+n}) g_{n+1}(X_{t+n+1}) \mid \mathcal{F}_t] = \mathbb{E}[g_0(X_t) \dots g_n(X_{t+n}) g_{n+1}(X_{t+n+1}) \mid X_t],$$

showing that  $(\mathcal{P}_{n+1})$  is true.

(d) Consider the vector space

$$\mathcal{H} = \{Y \in \sigma(X_s, s \geq t), \mathbb{E}[Y \mid \mathcal{F}_t] = \mathbb{E}[Y \mid X_t], \mathbb{P} - \text{a.s.}\}.$$

Let  $\{Y_n, n \in \mathbb{N}\}$  be an increasing sequence of nonnegative random variables in  $\mathcal{H}$  such that  $Y = \lim_{n \rightarrow \infty} Y_n$  is bounded. Show that

$$\mathbb{E}[Y \mid \mathcal{F}_t] = \lim_{n \rightarrow \infty} \mathbb{E}[Y_n \mid \mathcal{F}_t] = \lim_{n \rightarrow \infty} \mathbb{E}[Y_n \mid X_t] = \mathbb{E}[Y \mid X_t], \quad \mathbb{P} - \text{a.s.}$$

(e) Show that, for all  $t \in \mathbb{N}$  and all bounded  $\sigma(X_s, s \geq t)$ -measurable random variables  $Y$ ,

$$\mathbb{E}[Y \mid \mathcal{F}_t] = \mathbb{E}[Y \mid X_t], \quad \mathbb{P} - \text{a.s.} \quad (8.1)$$

**Exercise 8.3.** Let  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  be an adapted stochastic process such that for all  $t \in \mathbb{N}$  and all bounded  $\sigma(X_s, s \geq t)$ -measurable random variables  $Y$ ,  $\mathbb{E}[Y|\mathcal{F}_t] = \mathbb{E}[Y|X_t]$ ,  $\mathbb{P}$ -a.s.. For  $A \in \mathcal{F}_t$  and  $B \in \sigma(X_s, s \geq t)$ , denote  $Z = \mathbb{1}_A$  and  $Y = \mathbb{1}_B$ . Show that  $\mathbb{P}[A \cap B|\mathcal{F}_t] = Z\mathbb{E}[Y|X_t]$   $\mathbb{P}$ -a.s. and that  $\mathbb{P}[A \cap B|X_t] = \mathbb{P}[A|X_t]\mathbb{P}[B|X_t]$   $\mathbb{P}$ -a.s.

**Exercise 8.4.** Let  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  be an adapted stochastic process such that for all  $t \in \mathbb{N}$ ,  $A \in \mathcal{F}_t$  and  $B \in \sigma(X_s, s \geq t)$ ,  $\mathbb{P}[A \cap B|X_t] = \mathbb{P}[A|X_t]\mathbb{P}[B|X_t]$ ,  $\mathbb{P}$ -a.s.

- (a) Show that  $\mathbb{E}[YZ|X_t] = \mathbb{E}[Y|X_t]\mathbb{E}[Z|X_t]$  for all bounded  $\sigma(X_s, s \geq t)$ -measurable random variables  $Y$  and  $\mathcal{F}_t$ -measurable random variables  $Z$ .
- (b) Show that  $\{(X_t, \mathcal{F}_t), t \in \mathbb{N}\}$  is a Markov chain.

**Exercise 8.5.** Combine the following arguments to get the proof of Proposition 2.5.

- (a) Assume first that  $f$  is a simple nonnegative function, i.e.,  $f = \sum_{i \in I} \beta_i \mathbb{1}_{B_i}$  for a finite collection of nonnegative numbers  $\beta_i$  and sets  $B_i \in \mathcal{Y}$ . Then, using property (ii) of Theorem 2.2, for all  $x \in X$ ,  $Nf(x) = \sum_{i \in I} \beta_i N(x, B_i)$ , the function  $Nf$  is measurable.
- (b) Let now  $f \in \mathbb{F}_+(X)$  and let  $\{f_n, n \in \mathbb{N}\}$  be a nondecreasing sequence of measurable nonnegative simple functions such that  $\lim_{n \rightarrow \infty} f_n(x) = f(x)$ . Then, by the monotone convergence theorem, for all  $x \in X$ ,  $Nf(x) = \lim_{n \rightarrow \infty} Nf_n(x)$ .
- (c) Since any nonnegative measurable function is a pointwise limit of a nondecreasing sequence of nonnegative measurable functions, and since a pointwise limit of measurable function is measurable,  $Nf$  is measurable.

**Exercise 8.6.** Combine the following arguments to get the proof of Proposition 2.6.

- (a) It is easily seen that  $\mu N(A) \geq 0$  for all  $A \in \mathcal{Y}$  and  $\mu N(\emptyset) = 0$  ( $N(x, \emptyset) = 0$  for all  $x \in X$ ).
- (b) Let  $\{A_i, i \in \mathbb{N}\} \subset \mathcal{X}$  be a sequence of pairwise disjoint sets. For any integer  $n$ ,

$$\mu N\left(\bigcup_{i=1}^n A_i\right) = \int \mu(dx) N\left(x, \bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \int \mu(dx) N(x, A_i) = \sum_{i=1}^n \mu N(A_i).$$

The monotone convergence theorem therefore implies that

$$\mu N\left(\bigcup_{i=1}^{\infty} A_i\right) = \int \mu(dx) N\left(x, \bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \int \mu(dx) N(x, A_i) = \sum_{i=1}^{\infty} \mu N(A_i).$$

which establishes the  $\sigma$ -additivity.

**Exercise 8.7.** Prove Proposition 2.7.

**Exercise 8.8.** Prove Proposition 2.8.

**Exercise 8.9.** Let  $(X, \mathcal{X})$ ,  $(Y, \mathcal{Y})$ , and  $(Z, \mathcal{Z})$  be measurable spaces. Let  $M$  be a kernel on  $X \times \mathcal{Y}$  and  $N$  be a kernel on  $Y \times \mathcal{Z}$ .

- (a) Show that, if  $M$  and  $N$  are both finite (resp. bounded) kernels, then  $M \otimes N$  is finite (resp. bounded) kernel.
- (b) Show that, if  $M$  and  $N$  are both Markov kernels, then  $M \otimes N$  is a Markov kernel.
- (c) Show that, if  $(U, \mathcal{U})$  is a measurable space and  $P$  is a kernel on  $Z \times \mathcal{U}$ , then  $(M \otimes N) \otimes P = M \otimes (N \otimes P)$ , the tensor product of kernels is associative.

**Exercise 8.10 (Proof of Theorem 2.10: direct implication).** Denote by  $\mathcal{H}$  the set of measurable functions  $f \in \mathbb{F}_b(X^{t+1}, \mathcal{X}^{\otimes(t+1)})$  such that (2.6) is satisfied.

- (a) Show that  $\mathcal{H}$  is a vector space.
- (b) Let  $\{f_n, n \in \mathbb{N}\}$  be a nondecreasing sequence of nonnegative functions in  $\mathcal{H}$  such that  $\lim_{n \rightarrow \infty} f_n := f$  is bounded. Show that  $f$  belongs to  $\mathcal{H}$ .

(c) For  $t \geq 1$ , assume that (2.6) holds for  $t-1$  and  $f$  of the form

$$f_0 \otimes \cdots \otimes f_t(x_0, \dots, x_t) = f_0(x_0) \cdots f_t(x_t), \quad f_s \in \mathbb{F}_b(\mathbb{X}). \quad (8.2)$$

Show that

$$\begin{aligned} \mathbb{E} \left[ \prod_{s=0}^t f_s(X_s) \right] &= \mathbb{E} \left[ \prod_{s=0}^{t-1} f_s(X_s) \mathbb{E}[f_t(X_t) | \mathcal{F}_{t-1}] \right] = \mathbb{E} \left[ \prod_{s=0}^{t-1} f_s(X_s) P f_t(X_{t-1}) \right] \\ &= \nu \otimes P^{\otimes(t-1)}(f_0 \otimes \cdots \otimes f_{t-1} P f_t) = \nu \otimes P^{\otimes t}(f_0 \otimes \cdots \otimes f_t). \end{aligned}$$

(d) Show that  $\mathcal{H}$  contains the functions of the form (8.2).

(e) Conclude.

**Exercise 8.11 (Proof of Theorem 2.10: converse).** For  $t = 0$ , (2.6) implies that  $\nu$  is the law of  $X_0$ . We have to prove that, for all  $t \geq 1$ ,  $f \in \mathbb{F}_+(\mathbb{X})$  and  $\mathcal{F}_{t-1}^X$ -measurable random variable  $Y$ :

$$\mathbb{E}[f(X_t)Y] = \mathbb{E}[P f(X_{t-1})Y]. \quad (8.3)$$

Denote by  $\mathcal{H}$  the set of  $\mathcal{F}_{t-1}^X$ -measurable random variables  $Y$  satisfying (8.3).

(a) Show that  $\mathcal{H}$  is a vector space.

(b) Show that if  $\{Y_t, t \in \mathbb{N}\}$  is an increasing sequence of nonnegative random variables such that  $Y = \lim_{n \rightarrow \infty} Y_n$  is bounded, then  $Y \in \mathcal{H}$ .

(c) Show that (8.3) holds for  $Y = f_0(X_0)f_1(X_1) \cdots f_{t-1}(X_{t-1})$  where  $f_s \in \mathbb{F}_b(\mathbb{X})$  and conclude.

**Exercise 8.12 (Functional autoregressive process).** Consider the following recursion

$$X_t = g(X_{t-1}) + \sigma(X_{t-1})Z_t,$$

where  $\{Z_t, t \in \mathbb{N}\}$  is an i.i.d. sequence,  $g: \mathbb{R} \rightarrow \mathbb{R}$  and  $\sigma: \mathbb{R} \rightarrow \mathbb{R}_+^*$  are measurable functions.

(a) Determine the Markov kernel of this chain.

(b) Assume that the distribution of  $Z_0$  has a density  $q$  with respect to Lebesgue's measure on  $\mathbb{R}$ ; show that the Markov kernel  $P$  has a density. Give the expression of this density.

**Exercise 8.13 (Setwise convergence / existence of a stationary distribution).** Assume that for some initial probability  $\xi$ , the sequence of probabilities  $\{\xi^{P^t}, t \in \mathbb{N}\}$  converges setwise, i.e., for any  $A \in \mathcal{X}$ , the sequence  $\{\xi^{P^t}(A), t \in \mathbb{N}\}$  has a limit, denoted  $\gamma_\xi(A)$ , i.e.,  $\lim_{t \rightarrow \infty} \xi^{P^t}(A) = \gamma_\xi(A)$ .

(a) Show that,  $\gamma_\xi(A) = \gamma_\xi P(A)$ .

(b) Assume that there exists a unique invariant probability measure. Show that the limit  $\gamma_\xi$  is independent of the initial distribution  $\xi$ .

(c) Conversely, assume that for any initial distribution  $\xi$  on  $(\mathbb{X}, \mathcal{X})$ , the sequence  $\{\xi^{P^t}, t \in \mathbb{N}\}$  has the same setwise limit. Show that there exists a unique stationary distribution.

**Exercise 8.14 (Weak convergence / existence of a stationary distribution).** Let  $(\mathbb{X}, d)$  be a Polish space. We assume that for some initial measure  $\xi$ , the sequence of probability measures  $\{\xi^{P^t}, t \in \mathbb{N}\}$  converges weakly to  $\gamma_\xi$ . In addition, for any  $f \in C_b(\mathbb{X})$  (the set of bounded continuous functions),  $Pf \in C_b(\mathbb{X})$ .

(a) Show that  $\gamma_\xi = \gamma_\xi P$ .

(b) Show that if the kernel  $P$  has a unique invariant distribution, then the weak limit of the sequence  $\{\xi^{P^t}, t \in \mathbb{N}\}$  does not depend on the initial distribution  $\xi$ .

(c) Conversely, if for any initial distribution  $\xi$  the sequence of probability  $\{\xi^{P^t}, t \in \mathbb{N}\}$  converges weakly to the same limiting distribution, then the invariant distribution is unique.

**Exercise 8.15 (Existence and uniqueness of the stationary measure of an AR(1)).** Consider a first order autoregressive process  $X_t = \phi X_{t-1} + Z_t, t \geq 1$ , where  $\{Z_t, t \in \mathbb{N}\}$  is a zero mean i.i.d. sequence and  $X_0$  is independent of  $\{Z_t, t \in \mathbb{N}\}$  and is distributed according to some probability  $\xi$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . We assume that  $|\phi| < 1$ .

- (a) Show that  $X_t = \phi^t X_0 + \sum_{i=0}^{t-1} \phi^i Z_{t-i}$   
 (b) Show that, for each integer  $t$ , the random variable  $X_t$  has the same distribution as  $\phi^t X_0 + \sum_{i=0}^{t-1} \phi^i Z_{i+1}$ .  
 (c) Define  $\mathcal{F}_t^Z = \sigma(\{Z_s, 0 \leq s \leq t\})$  and

$$Y_t = \sum_{i=0}^{t-1} \phi^i Z_{i+1} = Y_{t-1} + \phi^{t-1} Z_t.$$

Show that  $\{(Y_t, \mathcal{F}_t^Z), t \in \mathbb{N}\}$  is an  $L^1$ -bounded martingale.

- (d) Show that the sequence  $\{Y_t, t \in \mathbb{N}\}$  converges  $\mathbb{P}$ -a.s. to the integrable random variable  $Y_\infty$ . We denote by  $\pi$  the law of  $Y_\infty$ .  
 (e) Show that for any function  $f$  and any initial distribution  $\xi$ ,  $\mathbb{E}[\xi][f(X_t)] \rightarrow \pi(f)$ .  
 (f) Using (8.14), show that  $\pi$  is the unique stationary distribution.

**Exercise 8.16.** The DAR(1) (see ??) is given by  $X_t = V_t X_{t-1} + (1 - V_t) Z_t$  where  $\{V_t\}$  is an i.i.d. sequence of binary random variables with  $\mathbb{P}(V_t = 1) = \alpha$ ,  $\{Z_t\}$  is an i.i.d. sequence of random variables with distribution given by  $\pi$  on  $(X, \mathcal{X})$  and  $\{V_t\}$  and  $\{Z_t\}$  are independent. The initial random variable  $X_0$  is assumed to be independent of  $\{V_t, t \in \mathbb{N}\}$  and  $\{Z_t, t \in \mathbb{N}\}$  and is distributed according to  $\xi$ .

- (a) Show that, for all  $t \geq 1$  and any bounded measurable function,

$$\mathbb{E}[\xi][f(X_t)] = \alpha \mathbb{E}[\xi][f(X_{t-1})] + (1 - \alpha) \pi(f).$$

- (b) Show that  $\mathbb{E}[\xi][f(X_{t-1})] = \pi(f)$ , then  $\mathbb{E}[\xi][f(X_t)] = \pi(f)$  and that  $\pi$  is the unique stationary distribution.  
 (c) Assume that  $X = \mathbb{N}$  and that  $\sum_{k=0}^{\infty} k^2 \pi(k) < \infty$ . For any positive integer  $h$ , show that  $\text{Cov}[\pi] X_h X_0 = \alpha \text{Cov}[\pi] X_{h-1} X_0$  and  $\text{Cov}[\pi] X_h X_0 = \alpha^h \text{Var}[\pi] X_0$ .  
 (d) Show that the DAR(1) process has exactly the same autocovariance structure as an AR(1) process.  
 (e) Explain why the DAR(1) process cannot exhibit negative dependence.

**Exercise 8.17.** Consider the self-excited threshold autoregressive model defined by

$$X_t = \phi_0^{(j)} + \sum_{i=1}^{p_j} \phi_i^{(j)} X_{t-i} + Z_t^{(j)} \quad \text{if } X_{t-d} \in (r_{j-1}, r_j] \quad (8.4)$$

where  $\{Z_t^{(j)}, t \in \mathbb{N}\} \sim_{\text{iid}} N(0, \sigma_j^2)$ , for  $j = 1, \dots, k$ , the positive integer  $d$  is a specified delay and  $r_0 = -\infty < r_1 < \dots < r_{p-1} < r_p = +\infty$  is a partition of  $\mathbb{R}$ .

- (a) Show that  $\{X_t, t \in \mathbb{N}\}$  is a  $m$ -th order Markov chain.  
 (b) Determine the Markov kernel of this chain.

**Exercise 8.18.** Let  $\{Z_t, t \in \mathbb{N}\}$  be Gaussian white noise with variance  $\sigma^2$  and let  $|\phi| < 1$  be a constant. Consider the process  $X_0 = Z_0$ , and  $X_t = \phi X_{t-1} + Z_t$ , for  $t = 1, 2, \dots$ .

- (a) Find the mean and the variance of  $X_t$ . Is  $\{X_t, t \in \mathbb{N}\}$  stationary?  
 (b) Show that for all  $h \in \{0, \dots, t\}$ ,

$$\text{Cor}(X_t, X_{t-h}) = \phi^h \left[ \frac{\text{Var}(X_{t-h})}{\text{Var}(X_t)} \right]^{1/2}.$$

- (c) Show that  $\lim_{t \rightarrow \infty} \text{Var}(X_t) = \sigma^2 / (1 - \phi^2)$  and that  $\lim_{t \rightarrow \infty} \text{Cor}(X_t, X_{t-h}) = \phi^h$ ,  $h \geq 0$ .  
 (d) Comment on how you could use these results to simulate  $n$  observations of a stationary Gaussian AR(1) model from simulated i.i.d.  $N(0, 1)$  values.  
 (e) Now suppose  $X_0 = Z_0 / \sqrt{1 - \phi^2}$ . Is  $\{X_t, t \in \mathbb{N}\}$  stationary?

**Exercise 8.19 (Simple Markovian bilinear model).** Consider the simple Markovian bilinear model

$$X_t = aX_{t-1} + bZ_t X_{t-1} + Z_t = (a + bZ_t) X_{t-1} + Z_t. \quad (8.5)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is an i.i.d. sequence of standard Gaussian variables and is independent of  $X_0 \sim \xi$ . We assume that  $a^2 + b^2 < 1$ .

(a) Show that, for all  $t \in \mathbb{N}$ ,

$$X_t = \prod_{i=1}^t (a + bZ_i) X_0 + Z_t + \sum_{j=1}^{t-1} Z_{t-j} \prod_{i=t-j+1}^t (a + bZ_i).$$

(b) Show that  $\lim_{t \rightarrow \infty} (\prod_{i=1}^t |a + bZ_i|)^{1/t}$  exists  $\mathbb{P}$ -a.s. and that this limit is strictly less than 1.

(c) Set  $Y_t = Z_0 + \sum_{j=1}^{t-1} Z_j \prod_{i=0}^{j-1} (a + bZ_i)$ . Show that the random variables  $Y_t$  and  $Z_t + \sum_{j=1}^{t-1} Z_{t-j} \prod_{i=t-j+1}^t (a + bZ_i)$  have the same distributions.

(d) Show that  $\{(Y_t, \mathcal{F}_t^Z), t \in \mathbb{N}\}$  is an  $L^2$ -bounded martingale converging  $\mathbb{P}$ -a.s. and in  $L^2(\mathbb{P})$  to  $Y_\infty = Z_0 + \sum_{j=1}^\infty Z_j \prod_{i=0}^{j-1} (a + bZ_i)$ .

(e) Show that for any  $f \in C_b(\mathbb{R})$ ,  $\lim_{t \rightarrow \infty} \mathbb{E}[\xi] [f(X_t)] = \pi(f)$ , where  $\pi_{a,b}$  is the distribution of  $Y_\infty$  (the subscripts  $(a, b)$  stress the dependence of the stationary distributions in  $a$  and  $b$ ).

(f) Using Exercise 8.14, show that  $\pi_{a,b}$  is the unique stationary distribution.

(g) Show that  $\int x \pi_{a,b}(dx) = 0$  and  $\int x^2 \pi_{a,b}(dx) = 1/(1 - a^2 - b^2)$ .

**Exercise 8.20 (Simple Markovian bilinear model cont.).** Consider again the simple bilinear model

$$X_t^{(a,b)} = aX_{t-1}^{(a,b)} + bZ_t X_{t-1}^{(a,b)} + Z_t = (a + bZ_t) X_{t-1}^{(a,b)} + Z_t, \quad (8.6)$$

where  $\{Z_t, t \in \mathbb{N}\}$  is an i.i.d. sequence of standard Gaussian variables and is independent of  $X_0 \sim \pi_{a,b}$  where  $\pi_{a,b}$  is the stationary distribution. The superscript  $(a, b)$  is used to stress the dependence of the process in these coefficients.

(a) Show that

$$\mathbb{E}|X_t^{(a,b)} - X_t^{(a,0)}|^2 = \frac{b^2}{(1 - a^2 - b^2)(1 - a^2)} \rightarrow 0, \quad \text{as } b \rightarrow 0.$$

(b) Assume that  $b > 0$ . Show that the cumulative distribution function  $x \mapsto F_{a,b}(x)$  of  $\pi_{a,b}$  is the solution of an integral equation.

(c) Show that, if  $x \neq a/b$ , the cumulative dist is differentiable at  $x$ .

(d) Show that  $F_{a,b}$  does not have a discrete component at  $x = -a/b$ .

(e) Conclude that the stationary distribution  $\pi_{a,b}$  has a density satisfying the following integral equation

$$f_{a,b}(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{f_{a,b}(s)}{|1 + bs|} \exp \left[ -\frac{1}{2} \left( \frac{x - as}{1 + bs} \right)^2 \right] \text{Leb}(ds).$$

(f) Show that  $f_{a,b}(-1/b) > 0$  and that  $\lim_{x \rightarrow -a/b} f_b(x) = \infty$ .

**Exercise 8.21.** Consider a GARCH( $p, q$ ) model

$$X_t = \sigma_t \varepsilon_t \quad (8.7)$$

$$\sigma_t^2 = \alpha_0 + \alpha_1 X_{t-1}^2 + \cdots + \alpha_p X_{t-p}^2 + \beta_1 \sigma_{t-1}^2 + \cdots + \beta_q \sigma_{t-q}^2, \quad (8.8)$$

where  $\{\varepsilon_t, t \in \mathbb{N}\}$  is an i.i.d. sequence with zero-mean and unit-variance, and the coefficients  $\alpha_j, j \in \{0, \dots, p\}$  and  $\beta_j, j \in \{1, \dots, q\}$  are nonnegative. State conditions upon which the GARCH( $p, q$ ) model admits a unique invariant stationary distribution.

**Exercise 8.22 (Basic Gibbs).** Suppose we wish to obtain samples from a bivariate normal distribution,

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \theta \\ \theta & 1 \end{pmatrix} \right],$$

where  $|\theta| < 1$ .

- (a) Determine the univariate conditionals  $p_1^\theta(X | Y)$  and  $p_2^\theta(Y | X)$ ; see (??)–(??).
- (b) Consider the Markov chain: Pick  $X^{(0)} = x_0$ , and then iterate the process  $X^{(0)} \mapsto Y^{(0)} \mapsto X^{(1)} \mapsto Y^{(1)} \mapsto \dots \mapsto X^{(t)} \mapsto Y^{(t)} \mapsto \dots$ , where  $Y^{(t)}$  is a sample from  $p_2^\theta(\cdot | X^{(t)})$ , and  $X^{(t)}$  is a sample from  $p_1^\theta(\cdot | Y^{(t-1)})$ . Write the joint distribution of  $(X^{(t)}, Y^{(t)})$  in terms of the starting value  $x_0$  and the correlation  $\theta$ .
- (c) What is the asymptotic distribution of  $(X^{(t)}, Y^{(t)})$  as  $t \rightarrow \infty$ ?
- (d) How can you use the results of part (c) to obtain a pseudo random sample of  $n$  bivariate normals using only pseudo samples from a univariate normal distribution? How does the value of  $\theta$  affect the procedure?

**Exercise 8.23 (Accept-reject algorithm).** We wish to sample the target density  $\pi$  (with respect to a dominating measure  $\lambda$ ), which is known up to a multiplicative constant. Let  $q$  be a proposal density (assumed to be easier to sample.) We assume that there exists a constant  $M < \infty$  such that  $\pi(x)/q(x) \leq M$  for all  $x \in X$ . The Accept-Reject goes as follows: first, we generate  $Y \sim g$  and, independently, we generate  $U \sim \text{Unif}([0, 1])$ . If  $Mq(Y)U \leq \pi(Y)$ , then we set  $X = Y$ . If the inequality is not satisfied, we then discard  $Y$  and  $U$  and start again.

- (a) Show that  $X$  is distributed according to  $\pi$ .
- (b) What is the distribution of the number of trials required per sample?
- (c) What is the average number of samples needed for one simulation?
- (d) Propose a method to estimate the normalizing constant of  $\pi$ .
- (e) Compare with the Metropolis-Hastings algorithm with Independent proposal.

**Exercise 8.24 (The independent sampler MCMC).** We wish to sample a target distribution  $\Gamma(\alpha, \beta)$  using the independent MCMC algorithm. We will use as the proposal distribution  $\gamma(\lfloor \alpha \rfloor, b)$ .

- (a) Derive the corresponding Accept-Reject method. Explain why, when  $\beta = 1$ , the choice of  $b$  minimizing the average number of simulations per sample is  $b = \lfloor \alpha \rfloor / \alpha$ .
- (b) Generate 10000  $\Gamma(5, 5/5.25)$  random variables to derive a  $\Gamma(5.25, 1)$  sample.
- (c) Use the same sample in the corresponding Metropolis-Hastings algorithm to generate 10000  $\Gamma(5.25, 1)$  random variables.
- (d) Compare the algorithms using (i) their acceptance rates and (ii) the estimates of the mean and variance of the  $\Gamma(5.25, 1)$  along with their errors.

**Exercise 8.25 (Functional autoregressive).** Let  $P$  be a Markov kernel on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . Show that  $(x, u) \mapsto G(x, u) = \inf\{x' \in \mathbb{R} : F(x, x') \geq u\}$  is a Borel function where  $F : (x, x') \mapsto P(x, (\infty, x'])$  is a Borel function.

**Exercise 8.26 (Reversibility of the Gibbs sampler).** We use the notations introduced in Section 4.1.5.

- (a) Show that for any functions  $f, g \in \mathbb{F}_+(\mathbf{X})$ ,

$$\iint f(x)g(x')\pi(x)\lambda(\mathrm{d}x)K_k(x, \mathrm{d}x') = \int \left\{ f(x)\pi(x)\lambda_k(\mathrm{d}x^{[k]}) \int g(x'^{[k]}, x^{[-k]})\pi_k(x'^{[k]} | x^{[-k]})\lambda_k(\mathrm{d}x'^{[k]}) \right\} \lambda_{-k}(\mathrm{d}x^{[-k]}),$$

where  $(x'^{[k]}, x^{[-k]})$  refers to the element  $u$  of  $\mathbf{X}$  such that  $u_k = x'^{[k]}$  and  $(u_{\ell \neq k}) = x^{[-k]}$

- (b) Show that

$$\pi(x^{[k]}, x^{[-k]})\pi_k(x'^{[k]} | x^{[-k]}) = \pi_k(x^{[k]} | x^{[-k]})\pi(x'^{[k]}, x^{[-k]}),$$

- (c) Conclude.

**Exercise 8.27 (A simple time series of counts).** Consider the following iterated random function:

$$X_{k+1} = d + aX_k + b[-X_k \ln(U_{k+1})], \quad k \geq 0 \quad (8.9)$$

where  $\{U_k, k \in \mathbb{N}^*\}$  is a sequence of i.i.d. random variables such that  $U_k \sim \text{Unif}[0, 1]$ . We assume that  $X_0$  is independent of  $\{U_k, k \in \mathbb{N}^*\}$  and that  $d, a, b > 0$  and  $a + b < 1$ . By simplification, we assume  $X_0 \geq d$ , so that  $\{X_k, k \in \mathbb{N}\}$  is a Markov chain taking values on  $X = [d, \infty)$ .

- (a) Show that for some function  $f$ , the recursion on the  $\{X_k, k \in \mathbb{N}\}$  may be written as: for all  $k \geq 0$ ,

$$\begin{aligned} Y_{k+1} | \mathcal{F}_k &\sim \text{Geom}(f(X_k)) , \\ X_{k+1} &= d + aX_k + b(Y_{k+1} - 1) . \end{aligned} \tag{8.10}$$

where  $\mathcal{F}_k = \sigma(X_0, \dots, X_k, Y_1, \dots, Y_k)$ .

- (b) Denote  $X_1(x) = d + ax + b \lfloor -x \ln(U) \rfloor$  where  $U \sim \text{Unif}[0, 1]$ . Compute  $\mathbb{E}[X_1(x)]$  explicitly.  
(c) Deduce that  $\varphi(x) : x \mapsto \mathbb{E}[X_1(x)]$  can be differentiated and show that  $|\varphi'(x)| < 1$  for all  $x \in \mathbb{X}$ .  
(d) By noting that  $X_1(x) \geq X_1(x')$  for all  $x \geq x'$ , deduce that there exists  $\rho < 1$  such that  $\mathbb{E}[|X_1(x) - X_1(x')|] \leq \rho|x - x'|$  for all  $x, x' \in \mathbb{X}$ . Conclude.





# Appendix A

## Technical results

### A.1 Limit theorems for triangular arrays

In this section, we derive limit theorems for triangular arrays of dependent random variables. Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space, let  $X$  be a random variable and let  $\mathcal{G}$  be a sub- $\sigma$  field of  $\mathcal{F}$ . Let  $\{M_N\}_{N \geq 0}$  be a sequence of positive integers,  $\{U_{N,i}\}_{i=1}^\infty$  be a triangular array of random variables, and  $\{\mathcal{F}_{N,i}\}_{0 \leq i \leq \infty}$  be a triangular array of sub-sigma-fields of  $\mathcal{F}$ . Throughout this section, it is assumed that  $\mathcal{F}_{N,i-1} \subseteq \mathcal{F}_{N,i}$  and for each  $N$  and  $i = 1, \dots, M_N$ ,  $U_{N,i}$  is  $\mathcal{F}_{N,i}$ -measurable. We preface the proof with some technical lemmas.

**Lemma A.1** *Let  $\{(Z_n, \mathcal{F}_n), n \in \mathbb{N}\}$  be an adapted sequence of nonnegative random variables. Then, for all  $\varepsilon > 0$ ,*

$$\xi_n = \mathbb{E} \left[ \sum_{i=1}^n Z_i \mathbb{1} \left( \sum_{j=1}^i \mathbb{E} [Z_j | \mathcal{F}_{j-1}] \right) \leq \varepsilon \right] \leq \varepsilon .$$

PROOF. We set  $\mu_i = \sum_{j=1}^i \mathbb{E} [Z_j | \mathcal{F}_{j-1}]$ . Then,

$$\xi_n = \sum_{i=1}^n \mathbb{E} [Z_i \mathbb{1} (\{\mu_i \leq \varepsilon\})] = \mathbb{E} \left[ \sum_{i=1}^n \mathbb{E} [Z_i | \mathcal{F}_{i-1}] \mathbb{1} (\{\mu_i \leq \varepsilon\}) \right] .$$

Set  $\tau := \max\{1 \leq i \leq n, \mu_i \leq \varepsilon\}$  and  $\tau = 0$  on  $\{\min 1 \leq i \leq n, \mu_i > \varepsilon\}$ . On  $\{\tau = 0\}$  we get  $\sum_{i=1}^n \mathbb{E} [Z_i | \mathcal{F}_{i-1}] \mathbb{1} (\{\mu_i \leq \varepsilon\}) = 0$  and, on  $\{\tau > 0\}$ , since  $\mu_1 \leq \mu_2 \leq \dots \leq \mu_n$ ,

$$\sum_{i=1}^n \mathbb{E} [Z_i | \mathcal{F}_{i-1}] \mathbb{1} (\{\mu_i \leq \varepsilon\}) = \sum_{i=1}^{\tau} \mathbb{E} [Z_i | \mathcal{F}_{i-1}] \leq \varepsilon .$$

■

The following Lemma has been established in [?], Lemma 3.5.

**Lemma A.2** *Let  $\{(Z_n, \mathcal{F}_n), n \in \mathbb{N}\}$  be an adapted sequence of nonnegative random variables. Then, for all  $\varepsilon > 0$ , and  $\alpha > 0$ ,*

$$\mathbb{P}(\max_{1 \leq i \leq n} Z_i > \varepsilon) \leq \alpha + \mathbb{P} \left( \sum_{i=1}^n \mathbb{P} [Z_i > \varepsilon | \mathcal{F}_{i-1}] > \alpha \right) .$$

PROOF. We set  $v_i := \sum_{j=1}^i \mathbb{P} [Z_j > \varepsilon | \mathcal{F}_{j-1}]$ . Then,

$$\mathbb{P}(\max_{1 \leq i \leq n} Z_i > \varepsilon) \leq \mathbb{P}(\max_{1 \leq i \leq n} Z_i > \varepsilon, v_n \leq \alpha) + \mathbb{P}(v_n > \alpha) ,$$

and

$$\begin{aligned} \mathbb{P}(\max_{1 \leq i \leq n} Z_i > \varepsilon, v_n \leq \alpha) &\leq \sum_{i=1}^n \mathbb{P}(Z_i > \varepsilon, v_n \leq \alpha) \\ &\leq \sum_{i=1}^n \mathbb{P}(Z_i > \varepsilon, v_i \leq \alpha) = \mathbb{E} \left[ \sum_{i=1}^n \mathbb{1}_{\{Z_i > \varepsilon\}} \mathbb{1}_{\{v_i \leq \alpha\}} \right] \leq \alpha , \end{aligned}$$

using Lemma A.1. ■

**Lemma A.3** *Let  $\mathcal{G}$  be a  $\sigma$ -field and  $X$  a random variable such that  $\mathbb{E} [X^2 | \mathcal{G}] < \infty$ . Then, for any  $\varepsilon > 0$ ,*

$$4\mathbb{E} [|X|^2 \mathbb{1}_{\{|X| \geq \varepsilon\}} | \mathcal{G}] \geq \mathbb{E} [|X - \mathbb{E} [X | \mathcal{G}]|^2 \mathbb{1}_{\{|X - \mathbb{E} [X | \mathcal{G}]| \geq 2\varepsilon\}} | \mathcal{G}] .$$

PROOF. Let  $Y = X - \mathbb{E} [X | \mathcal{G}]$ . We have  $\mathbb{E} [Y | \mathcal{G}] = 0$ . It is equivalent to show that for any  $\mathcal{G}$ -measurable random variable  $Z$ ,

$$\mathbb{E} [Y^2 \mathbb{1}_{\{|Y| \geq 2\varepsilon\}} | \mathcal{G}] \leq 4\mathbb{E} [|Y + Z|^2 \mathbb{1}_{\{|Y + Z| \geq \varepsilon\}} | \mathcal{G}] .$$

On the set  $\{|Z| < \varepsilon\}$ ,

$$\begin{aligned}
\mathbb{E} [Y^2 \mathbb{1}\{|Y| \geq 2\varepsilon\} | \mathcal{G}] &\leq 2\mathbb{E} [(Y+Z)^2 + Z^2] \mathbb{1}\{|Y+Z| \geq \varepsilon\} | \mathcal{G}] \\
&\leq 2(1+Z^2/\varepsilon^2) \mathbb{E} [(Y+Z)^2 \mathbb{1}\{|Y+Z| \geq \varepsilon\} | \mathcal{G}] \\
&\leq 4\mathbb{E} [(Y+Z)^2 \mathbb{1}\{|Y+Z| \geq \varepsilon\} | \mathcal{G}].
\end{aligned}$$

Moreover, on the set  $\{|Z| \geq \varepsilon\}$ , using that  $\mathbb{E}[ZY | \mathcal{G}] = Z\mathbb{E}[Y | \mathcal{G}] = 0$ .

$$\begin{aligned}
\mathbb{E} [Y^2 \mathbb{1}\{|Y| \geq 2\varepsilon\} | \mathcal{G}] &\leq \mathbb{E} [Y^2 + Z^2 - \varepsilon^2 | \mathcal{G}] \\
&\leq \mathbb{E} [(Y+Z)^2 - \varepsilon^2 | \mathcal{G}] \leq \mathbb{E} [(Y+Z)^2 \mathbb{1}\{|Y+Z| \geq \varepsilon\} | \mathcal{G}].
\end{aligned}$$

The proof is completed.  $\blacksquare$

**Theorem A.4.** Assume that  $\mathbb{E}[|U_{N,j}| | \mathcal{F}_{N,j-1}] < \infty$   $\mathbb{P}$ -a.s. for any  $N$  and any  $j = 1, \dots, M_N$ , and

$$\sup_N \mathbb{P} \left( \sum_{j=1}^{M_N} \mathbb{E}[|U_{N,j}| | \mathcal{F}_{N,j-1}] \geq \lambda \right) \rightarrow 0 \quad \text{as } \lambda \rightarrow \infty \quad (\text{A.1})$$

$$\sum_{j=1}^{M_N} \mathbb{E}[|U_{N,j}| \mathbb{1}\{|U_{N,j}| \geq \varepsilon\} | \mathcal{F}_{N,j-1}] \xrightarrow{\mathbb{P}\text{-prob}} 0 \quad \text{for any } \varepsilon > 0. \quad (\text{A.2})$$

Then,  $\max_{1 \leq i \leq M_N} \left| \sum_{j=1}^i U_{N,j} - \sum_{j=1}^i \mathbb{E}[U_{N,j} | \mathcal{F}_{N,j-1}] \right| \xrightarrow{\mathbb{P}\text{-prob}} 0$ .

PROOF. Assume first that for each  $N$  and each  $i = 1, \dots, M_N$ ,  $U_{N,i} \geq 0$ ,  $\mathbb{P}$ -a.s. By Lemma A.2, we have that for any constants  $\varepsilon$  and  $\eta > 0$ ,

$$\mathbb{P} \left[ \max_{1 \leq i \leq M_N} U_{N,i} \geq \varepsilon \right] \leq \eta + \mathbb{P} \left[ \sum_{i=1}^{M_N} \mathbb{P}(U_{N,i} \geq \varepsilon | \mathcal{F}_{N,i-1}) \geq \eta \right].$$

From the conditional version of the Chebyshev identity,

$$\mathbb{P} \left[ \max_{1 \leq i \leq M_N} U_{N,i} \geq \varepsilon \right] \leq \eta + \mathbb{P} \left[ \sum_{i=1}^{M_N} \mathbb{E}[U_{N,i} \mathbb{1}\{U_{N,i} \geq \varepsilon\} | \mathcal{F}_{N,i-1}] \geq \eta \varepsilon \right]. \quad (\text{A.3})$$

Let  $\varepsilon$  and  $\lambda > 0$  and define

$$\bar{U}_{N,i} := U_{N,i} \mathbb{1}\{U_{N,i} < \varepsilon\} \mathbb{1} \left\{ \sum_{j=1}^i \mathbb{E}[U_{N,j} | \mathcal{F}_{N,j-1}] < \lambda \right\}.$$

For any  $\delta > 0$ ,

$$\begin{aligned}
&\mathbb{P} \left( \max_{1 \leq i \leq M_N} \left| \sum_{j=1}^i U_{N,j} - \sum_{j=1}^i \mathbb{E}[U_{N,j} | \mathcal{F}_{N,j-1}] \right| \geq 2\delta \right) \\
&\leq \mathbb{P} \left( \max_{1 \leq i \leq M_N} \left| \sum_{j=1}^i \bar{U}_{N,j} - \sum_{j=1}^i \mathbb{E}[\bar{U}_{N,j} | \mathcal{F}_{N,j-1}] \right| \geq \delta \right) \\
&\quad + \mathbb{P} \left( \max_{1 \leq i \leq M_N} \left| \sum_{j=1}^i U_{N,j} - \bar{U}_{N,j} - \sum_{j=1}^i \mathbb{E}[U_{N,j} - \bar{U}_{N,j} | \mathcal{F}_{N,j-1}] \right| \geq \delta \right).
\end{aligned}$$

The second term in the right-hand side is bounded by

$$\begin{aligned}
&\mathbb{P} \left( \max_{1 \leq i \leq M_N} U_{N,i} \geq \varepsilon \right) + \mathbb{P} \left( \sum_{j=1}^{M_N} \mathbb{E}[U_{N,j} | \mathcal{F}_{N,j-1}] \geq \lambda \right) + \\
&\quad \mathbb{P} \left( \sum_{j=1}^{M_N} \mathbb{E}[U_{N,j} \mathbb{1}\{U_{N,j} \geq \varepsilon\} | \mathcal{F}_{N,j-1}] \geq \delta \right).
\end{aligned}$$

Eqs. (A.2) and (A.3) imply that the first and last terms in the last expression converge to zero for any  $\varepsilon > 0$  and (A.1) implies that the second term may be arbitrarily small by choosing for  $\lambda$  sufficiently large. Now, by the Doob maximal inequality,

$$\begin{aligned} \mathbb{P} \left( \max_{1 \leq i \leq M_N} \left| \sum_{j=1}^i \bar{U}_{N,j} - \mathbb{E} [\bar{U}_{N,j} | \mathcal{F}_{N,j-1}] \right| \geq \delta \right) \\ \leq \delta^{-2} \mathbb{E} \left[ \sum_{j=1}^{M_N} \mathbb{E} \left[ (\bar{U}_{N,j} - \mathbb{E} [\bar{U}_{N,j} | \mathcal{F}_{N,j-1}])^2 \middle| \mathcal{F}_{N,0} \right] \right]. \end{aligned}$$

This last term does not exceed

$$\begin{aligned} \delta^{-2} \mathbb{E} \left[ \sum_{i=1}^{M_N} \mathbb{E} [\bar{U}_{N,i}^2 | \mathcal{F}_{N,0}] \right] &\leq \delta^{-2} \varepsilon \mathbb{E} \left[ \sum_{j=1}^{M_N} \mathbb{E} [\bar{U}_{N,j} | \mathcal{F}_{N,0}] \right] \\ &\leq \delta^{-2} \varepsilon \mathbb{E} \left[ \sum_{j=1}^{M_N} \mathbb{E} [\bar{U}_{N,j} | \mathcal{F}_{N,j-1}] \right] \leq \delta^{-2} \varepsilon \lambda. \end{aligned}$$

Since  $\varepsilon$  is arbitrary, the proof follows for  $U_{N,j} \geq 0$ ,  $\mathbb{P}$ -a.s., for each  $N$  and  $j = 1, \dots, M_N$ . The proof extends to an arbitrary triangular array  $\{U_{N,j}\}_{i=1}^{M_N}$  by applying the preceding result to  $\{U_{N,j}^+\}_{1 \leq j \leq M_N}$  and  $\{U_{N,j}^-\}_{1 \leq j \leq M_N}$ . ■

**Lemma A.5** Assume that for all  $N$ ,  $\sum_{i=1}^{M_N} \mathbb{E} [U_{N,i}^2 | \mathcal{F}_{N,i-1}] = 1$ ,  $\mathbb{E} [U_{N,i} | \mathcal{F}_{N,i-1}] = 0$  for  $i = 1, \dots, M_N$ , and for all  $\varepsilon > 0$ ,

$$\sum_{i=1}^{M_N} \mathbb{E} [U_{N,i}^2 \mathbb{1}\{|U_{N,i}| \geq \varepsilon\} | \mathcal{F}_{N,0}] \xrightarrow{\mathbb{P}\text{-prob}} 0. \quad (\text{A.4})$$

Then, for any real  $u$ ,  $\mathbb{E} \left[ \exp \left( iu \sum_{j=1}^{M_N} U_{N,j} \right) \middle| \mathcal{F}_{N,0} \right] - \exp(-u^2/2) \xrightarrow{\mathbb{P}\text{-prob}} 0$ .

PROOF. Denote  $\sigma_{N,i}^2 := \mathbb{E} [U_{N,i}^2 | \mathcal{F}_{N,i-1}]$ . Write the following decomposition (with the convention  $\sum_{j=a}^b = 0$  if  $a > b$ ):

$$e^{iu \sum_{j=1}^{M_N} U_{N,j}} - e^{-\frac{u^2}{2} \sum_{j=1}^{M_N} \sigma_{N,j}^2} = \sum_{l=1}^{M_N} e^{iu \sum_{j=1}^{l-1} U_{N,j}} \left( e^{iu U_{N,l}} - e^{-\frac{u^2}{2} \sigma_{N,l}^2} \right) e^{-\frac{u^2}{2} \sum_{j=l+1}^{M_N} \sigma_{N,j}^2}.$$

Since  $\sum_{j=1}^{l-1} U_{N,j}$  and  $\sum_{j=l+1}^{M_N} \sigma_{N,j}^2 = 1 - \sum_{j=1}^l \sigma_{N,j}^2$  are  $\mathcal{F}_{N,l-1}$ -measurable,

$$\begin{aligned} \left| \mathbb{E} \left[ \exp \left( iu \sum_{j=1}^{M_N} U_{N,j} \right) - \exp \left( -(u^2/2) \sum_{j=1}^{M_N} \sigma_{N,j}^2 \right) \middle| \mathcal{F}_{N,0} \right] \right| \\ \leq \sum_{l=1}^{M_N} \mathbb{E} \left[ \left| \mathbb{E} \left[ \exp(iu U_{N,l}) \middle| \mathcal{F}_{N,l-1} \right] - \exp(-u^2 \sigma_{N,l}^2/2) \right| \middle| \mathcal{F}_{N,0} \right]. \quad (\text{A.5}) \end{aligned}$$

For any  $\varepsilon > 0$ , it is easily shown that

$$\begin{aligned} \mathbb{E} \left[ \sum_{l=1}^{M_N} \left| \mathbb{E} \left[ \exp(iu U_{N,l}) - 1 + \frac{1}{2} u^2 \sigma_{N,l}^2 \middle| \mathcal{F}_{N,l-1} \right] \right| \middle| \mathcal{F}_{N,0} \right] \\ \leq \frac{1}{6} \varepsilon |u|^3 + u^2 \sum_{l=1}^{M_N} \mathbb{E} [U_{N,l}^2 \mathbb{1}\{|U_{N,l}| \geq \varepsilon\} | \mathcal{F}_{N,0}] . \quad (\text{A.6}) \end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary, it follows from (A.4) that the right-hand side tends in probability to 0 as  $N \rightarrow \infty$ . Finally, for all  $\varepsilon > 0$ ,

$$\begin{aligned} \mathbb{E} \left[ \sum_{l=1}^{M_N} \left| \mathbb{E} \left[ \exp(-u^2 \sigma_{N,l}^2/2) - 1 + \frac{1}{2} u^2 \sigma_{N,l}^2 \middle| \mathcal{F}_{N,l-1} \right] \right| \middle| \mathcal{F}_{N,0} \right] \\ \leq \frac{u^4}{8} \sum_{l=1}^{M_N} \mathbb{E} [\sigma_{N,l}^4 | \mathcal{F}_{N,0}] \leq \frac{u^4}{8} \left( \varepsilon^2 + \sum_{j=1}^{M_N} \mathbb{E} [U_{N,j}^2 \mathbb{1}\{|U_{N,j}| \geq \varepsilon\} | \mathcal{F}_{N,0}] \right). \end{aligned}$$

(A.4) shows that the right-hand side of previous equation tends in probability to 0 as  $N \rightarrow \infty$ . The proof follows. ■

**Theorem A.6.** Assume that for each  $N$  and  $i = 1, \dots, M_N$ ,  $\mathbb{E} [U_{N,i}^2 | \mathcal{F}_{N,i-1}] < \infty$  and

$$\sum_{i=1}^{M_N} \{\mathbb{E}[U_{N,i}^2 | \mathcal{F}_{N,i-1}] - (\mathbb{E}[U_{N,i} | \mathcal{F}_{N,i-1}])^2\} \xrightarrow{\mathbb{P}\text{-prob}} \sigma^2 \quad \text{for some } \sigma^2 > 0, \quad (\text{A.7})$$

$$\sum_{i=1}^{M_N} \mathbb{E}[U_{N,i}^2 \mathbb{1}_{\{|U_{N,i}| \geq \varepsilon\}} | \mathcal{F}_{N,i-1}] \xrightarrow{\mathbb{P}\text{-prob}} 0 \quad \text{for any } \varepsilon > 0. \quad (\text{A.8})$$

Then, for any real  $u$ ,

$$\mathbb{E} \left[ \exp \left( iu \sum_{i=1}^{M_N} \{U_{N,i} - \mathbb{E}[U_{N,i} | \mathcal{F}_{N,i-1}]\} \right) \middle| \mathcal{F}_{N,0} \right] \xrightarrow{\mathbb{P}\text{-prob}} \exp(-(u^2/2)\sigma^2). \quad (\text{A.9})$$

PROOF. We first assume that  $\mathbb{E}[U_{N,i} | \mathcal{F}_{N,i-1}] = 0$  for all  $i = 1, \dots, M_N$ , and  $\sigma^2 = 1$ . Define  $\tau_N := \max \{1 \leq k \leq M_N : \sum_{j=1}^k \sigma_{N,j}^2 \leq 1\}$ , with the convention  $\max \emptyset = 0$ . Put  $\bar{U}_{N,k} = U_{N,k}$  for  $k \leq \tau_N$ ,  $\bar{U}_{N,k} = 0$  for  $\tau_N < k \leq M_N$  and  $\bar{U}_{N,M_N+1} = \left(1 - \sum_{j=1}^{\tau_N} \sigma_{N,j}^2\right)^{1/2} Y_N$ , where  $\{Y_N\}$  are  $\mathcal{N}(0, 1)$  independent and independent of  $\mathcal{F}_{N,M_N}$ . Put

$$\sum_{j=1}^{M_N} U_{N,j} = \sum_{j=1}^{M_N+1} \bar{U}_{N,j} - \bar{U}_{N,M_N+1} + \sum_{j=\tau_N+1}^{M_N} U_{N,j}. \quad (\text{A.10})$$

We will prove that (a)  $\{\bar{U}_{N,j}\}_{1 \leq j \leq M_N+1}$  satisfies the assumptions of Lemma A.5, (b)  $\bar{U}_{N,M_N+1} \xrightarrow{\mathbb{P}\text{-prob}} 0$ , (c)  $\sum_{j=\tau_N+1}^{M_N} U_{N,j} \xrightarrow{\mathbb{P}\text{-prob}} 0$ . If  $\tau_N < M_N$ , then for any  $\varepsilon > 0$ ,

$$0 \leq 1 - \sum_{j=1}^{\tau_N} \sigma_{N,j}^2 \leq \sigma_{N,\tau_N+1}^2 \leq \max_{1 \leq j \leq M_N} \sigma_{N,j}^2 \leq \varepsilon^2 + \sum_{j=1}^{M_N} \mathbb{E}[U_{N,j}^2 \mathbb{1}_{\{|U_{N,j}| \geq \varepsilon\}} | \mathcal{F}_{N,j-1}]$$

Since  $\varepsilon > 0$  is arbitrary, it follows from (A.8) that  $1 - \sum_{j=1}^{\tau_N} \sigma_{N,j}^2 \xrightarrow{\mathbb{P}\text{-prob}} 0$ , which implies that  $\mathbb{E}[\bar{U}_{N,M_N+1}^2 | \mathcal{F}_{N,0}] \xrightarrow{\mathbb{P}\text{-prob}} 0$ , showing (a) and (b). It remains to prove (c). We have

$$\sum_{j=\tau_N+1}^{M_N} \sigma_{N,j}^2 = \sum_{j=1}^{M_N} \sigma_{N,j}^2 - 1 + \left(1 - \sum_{j=1}^{\tau_N} \sigma_{N,j}^2\right) \xrightarrow{\mathbb{P}\text{-prob}} 0. \quad (\text{A.11})$$

For any  $\lambda > 0$ ,

$$\left( \mathbb{E} \left[ \sum_{j=\tau_N+1}^{M_N} U_{N,j} \mathbb{1}_{\left\{ \sum_{i=\tau_N+1}^j \sigma_{N,i}^2 \leq \lambda \right\}} \right] \right)^2 = \mathbb{E} \left[ \sum_{j=\tau_N+1}^{M_N} \sigma_{N,j}^2 \mathbb{1}_{\left\{ \sum_{i=\tau_N+1}^j \sigma_{N,i}^2 \leq \lambda \right\}} \right].$$

The term between braces converges to 0 in probability by (A.11) and its value is bounded by  $\lambda$ , which shows that  $\sum_{j=\tau_N+1}^{M_N} U_{N,j} \mathbb{1}_{\left\{ \sum_{i=\tau_N+1}^j \sigma_{N,i}^2 \leq \lambda \right\}} \xrightarrow{\mathbb{P}\text{-prob}} 0$ . Moreover,

$$\mathbb{P} \left( \sum_{j=\tau_N+1}^{M_N} U_{N,j} \mathbb{1}_{\left\{ \sum_{i=\tau_N+1}^j \sigma_{N,i}^2 > \lambda \right\}} \neq 0 \right) \leq \mathbb{P} \left( \sum_{i=\tau_N+1}^{M_N} \sigma_{N,i}^2 > \lambda \right),$$

which converges to 0 by (A.11). The proof is completed when  $\mathbb{E}[U_{N,i} | \mathcal{F}_{N,i-1}] = 0$ . To deal with the general case, it suffices to set  $\bar{U}_{N,i} = U_{N,i} - \mathbb{E}[U_{N,i} | \mathcal{F}_{N,i-1}]$  and use Lemma A.3.  $\blacksquare$

We may now specialize these results to stationary martingale increment sequences. This result was established by ?.

**Theorem A.7.** Assume that  $\{X_k, k \in \mathbb{N}\}$  is a strict-sense stationary, ergodic process such that  $\mathbb{E}[X_1^2]$  is finite and  $\mathbb{E}[X_k | \mathcal{F}_{k-1}^X] = 0$ , where  $\{\mathcal{F}_k^X, k \in \mathbb{N}\}$  is the natural filtration. Then,  $n^{-1/2} \sum_{k=1}^n X_k \xrightarrow{\mathcal{L}_P} \mathcal{N}(0, \mathbb{E}[X_1^2])$ .



