

Warsaw University  
Faculty of Mathematics, Informatics and Mechanics

Krzysztof Łatuszyński

Regeneration and Fixed-Width Analysis of  
Markov Chain Monte Carlo Algorithms

*PhD dissertation*

Supervisor  
dr hab. Wojciech Niemiro

Institute of Applied Mathematics and Mechanics  
Warsaw University

February 2008

Author's declaration:

aware of legal responsibility I hereby declare that I have written this dissertation myself and all the contents of the dissertation have been obtained by legal means.

February 18, 2008

*date*

.....

*Krzysztof Łatuszyński*

Supervisor's declaration:

the dissertation is ready to be reviewed

February 18, 2008

*date*

.....

*dr hab. Wojciech Niemiro*

## Abstract

In the thesis we take the split chain approach to analyzing Markov chains and use it to establish fixed-width results for estimators obtained via Markov chain Monte Carlo procedures (MCMC). Theoretical results include necessary and sufficient conditions in terms of regeneration for central limit theorems for ergodic Markov chains and a regenerative proof of a CLT version for uniformly ergodic Markov chains with  $E_\pi f^2 < \infty$ . To obtain asymptotic confidence intervals for MCMC estimators, strongly consistent estimators of the asymptotic variance are essential. We relax assumptions required to obtain such estimators. Moreover, under a drift condition, nonasymptotic fixed-width results for MCMC estimators for a general state space setting (not necessarily compact) and not necessarily bounded target function  $f$  are obtained. The last chapter is devoted to the idea of adaptive Monte Carlo simulation and provides convergence results and law of large numbers for adaptive procedures under path-stability condition for transition kernels.

**Keywords and phrases:** Markov chain, MCMC, adaptive Monte Carlo, split chain, regeneration, drift condition,  $(\varepsilon-\alpha)$ -approximation, confidence intervals, asymptotic confidence intervals, central limit theorem, law of large numbers

**AMS Subject Classification:** 60J10, 60J05, 60F15, 60F05

## Streszczenie

W pracy przedstawione są rezultaty dotyczące estymacji stałoprecyzyjnej dla algorytmów Monte Carlo opartych na łańcuchach Markowa (MCMC). Podstawową techniką w analizie łańcuchów Markowa i związanych z nimi procedur MCMC, jest łańcuch rozszczepiony i regeneracja, co prowadzi do koniecznego i dostatecznego warunku w terminach regeneracji dla centralnego twierdzenia granicznego dla ergodycznych łańcuchów Markowa. Dodatkowym rezultatem jest regeneracyjny dowód CTG dla jednostajnie ergodycznych łańcuchów Markowa przy założeniu  $E_\pi f^2 < \infty$ . Aby otrzymać asymptotyczne przedziały ufności za pomocą algorytmów MCMC konieczna jest m.in. mocno zgodna estymacja wariancji asymptotycznej. Osłabiamy znane założenia wymagane do konstrukcji takich estymatorów. Przy założeniu warunku dryfu, ale bez założeń o ograniczoności funkcji podcałkowej  $f$  i zwartości przestrzeni stanów, otrzymujemy nieasymptotyczną estymację stałoprecyzyjną. Ostatni rozdział poświęcony jest procedurom adaptacyjnym, a uzyskane tam wyniki dotyczące zbieżności i prawa wielkich liczb zakładają stabilność operatorów przejścia względem trajektorii.

**Słowa kluczowe:** łańcuch Markowa, MCMC, adaptacyjne Monte Carlo, łańcuch rozszczepiony, regeneracja, warunek dryfu,  $(\varepsilon - \alpha)$ -aproxymacja, przedziały ufności, asymptotyczne przedziały ufności, centralne twierdzenie graniczne, prawo wielkich liczb

**Klasyfikacja tematyczna wg. AMS:** 60J10, 60J05, 60F15, 60F05

# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	Markov Chain Monte Carlo . . . . .	7
1.2	Sampling Schemes and MCMC Algorithms . . . . .	8
1.2.1	The Metropolis Algorithm . . . . .	11
1.2.2	The Gibbs Sampler . . . . .	12
1.3	Overview of the Results . . . . .	13
<b>2</b>	<b>Some Markov Chains</b>	<b>17</b>
2.1	Stationarity and Ergodicity . . . . .	17
2.2	Small Sets and the Split Chain . . . . .	22
<b>3</b>	<b>A Complete Characterisation of <math>\sqrt{n}</math>-CLTs for Ergodic Markov Chains via Regeneration</b>	<b>28</b>
3.1	CLTs for Markov Chains . . . . .	28
3.2	Tools and Preliminary Results . . . . .	30
3.3	A Characterization of $\sqrt{n}$ -CLTs . . . . .	35
3.4	Uniform Ergodicity . . . . .	39
3.5	The difference between $m = 1$ and $m \neq 1$ . . . . .	43
<b>4</b>	<b>Fixed-Width Asymptotics</b>	<b>44</b>
4.1	Asymptotic Confidence Intervals . . . . .	44
4.2	Estimating Asymptotic Variance . . . . .	46
4.2.1	Batch Means . . . . .	46
4.2.2	Regenerative Estimation . . . . .	47
4.3	A Lemma and its Consequences . . . . .	49
<b>5</b>	<b>Fixed-Width Nonasymptotic Results under Drift Condition</b>	<b>53</b>
5.1	Introduction . . . . .	53

5.2	A Drift Condition and Preliminary Lemmas . . . . .	55
5.3	MSE Bounds . . . . .	59
5.4	$(\varepsilon - \alpha)$ -Approximation . . . . .	62
5.5	A Toy Example - Contracting Normals . . . . .	66
5.6	The Example - a Hierarchical Random Effects Model . . . . .	67
	5.6.1 The Model . . . . .	68
	5.6.2 Gibbs Samplers for the Model . . . . .	69
	5.6.3 Relations between Drift Conditions . . . . .	71
	5.6.4 Drift and Minorization Conditions for the Samplers . . . . .	73
	5.6.5 Obtaining the Bounds . . . . .	79
5.7	Concluding Remarks . . . . .	79
5.8	Appendix - Formulas for $\rho$ and $M$ . . . . .	81
	5.8.1 Formulas for general operators . . . . .	81
	5.8.2 Formulas for self-adjoint operators . . . . .	82
	5.8.3 Formulas for self-adjoint positive operators . . . . .	83
<b>6</b>	<b>Convergence Results for Adaptive Monte Carlo</b>	<b>84</b>
6.1	Introduction . . . . .	84
6.2	One Intuitive and One Not-so-Intuitive Example . . . . .	85
6.3	Convergence Results . . . . .	87
6.4	Proofs . . . . .	90
6.5	Appendix - Mixingales . . . . .	97
	<b>Bibliography</b>	<b>98</b>

# Chapter 1

## Introduction

In this chapter we give some background for results presented in later chapters and introduce main ideas behind the thesis in an informal way. Therefore mathematical rigour will not always be our priority here. We start with defining the problem addressed by Markov chain Monte Carlo methods in Section 1.1 and proceed to describing typical sampling schemes and MCMC algorithms (the Metropolis algorithm and the Gibbs sampler) in Section 1.2. Section 1.3 provides an overview of the results of the thesis.

### 1.1 Markov Chain Monte Carlo

Let  $\mathcal{X}$  be a region in a possibly high-dimensional space, and let  $f$  be a real valued function on  $\mathcal{X}$ . Moreover consider a probability distribution  $\pi$  with density  $p$  with respect to some standard measure  $dx$ , usually either Lebesgue or counting measure, i.e.  $\pi(dx) = p(x)dx$ . An essential part of many problems in Bayesian inference, statistical physics and combinatorial enumeration is the computation of analytically intractable integral

$$I = E_{\pi}f = \int_{\mathcal{X}} f(x)\pi(dx), \quad (1.1)$$

where  $p$  and thus  $\pi$  is known up to a normalizing constant and direct simulation from  $\pi$  is not feasible (see e.g. [Casella & Robert 1999], [Liu, JS 2001]). The common approach to this problem is to simulate an ergodic Markov chain  $(X_n)_{n \geq 0}$ , using a transition kernel  $P$ , with stationary distribution  $\pi$ , which ensures the convergence in distribution of  $X_n$  to a random variable from  $\pi$ .

Thus, for a "large enough"  $t$ ,  $X_n$  for  $n \geq t$  can be considered as having distribution approximately equal to  $\pi$ . Since a simple and powerful algorithm for constructing such a Markov chain has been introduced in 1953 by Metropolis et al. in the very seminal paper [Metropolis et al. 1953], various sampling schemes and approximation strategies for estimating the unknown value of  $I$  have been developed and analyzed ([Niemi & Pokarowski 2007], [Liu, JS 2001], [Casella & Robert 1999]). The method is referred to as Markov chain Monte Carlo (MCMC).

To avoid problems with integrating functions with respect to probability distributions with unknown normalizing constants, Bayesian statisticians used to restrict attention to conjugate priors (see e.g. [Robert 1994]). This concept, although technically appealing, deprives the Bayesian approach of flexibility which is one of its main strengths. Also, when building complex models with many parameters (as in the example of Section ??), even using conjugate priors usually leads to intractable multidimensional posterior distributions.

The invention of MCMC has transformed dramatically Bayesian inference since it allows practitioners to sample from complicated posterior distributions and to integrate functions with respect to these distributions. Thus Bayesian inference became a feasible and powerful approach for practitioners and now receives immense attention from the statistics community ([Roberts & Rosenthal 2005],[Casella & Robert 1999]).

In addition to their importance for applications, MCMC algorithms raise numerous questions related to Markov chains and probability. It is crucial to understand the nature and speed of convergence of the distribution of  $X_n$  to  $\pi$  as  $n \rightarrow \infty$ .

## 1.2 Sampling Schemes and MCMC Algorithms

Before we proceed to the description of MCMC algorithms let us recall the independent Monte Carlo solution to the problem in (1.1) when simulating from  $\pi$  is feasible. In this case one takes i.i.d. random variables  $X_1, \dots, X_n \sim \pi$  and estimates  $I$  by

$$\hat{I}_n = \frac{1}{n} \sum_{i=1}^n f(X_i). \quad (1.2)$$

*Remark 1.2.1.* Basic properties of the independent Monte Carlo estimation



are very easy to obtain.

- If  $I$  exists then  $\hat{I}_n$  is its unbiased and (by the weak law of large numbers) consistent estimate.
- Furthermore, if  $\pi f^2 < \infty$ , then by the classical Central Limit Theorem

$$\sqrt{n}(\hat{I} - I) \xrightarrow{d} N(0, \pi f^2 - (\pi f)^2).$$

- Confidence intervals for  $I$  can be obtained e.g. by the Chebyshev inequality

$$P(|\hat{I}_n - I| \geq \varepsilon) \leq \frac{\pi f^2 - (\pi f)^2}{n\varepsilon^2},$$

provided that the variance  $\pi f^2 - (\pi f)^2$  can be bounded a priori.

- Asymptotic confidence intervals can be derived from the CLT,

$$P(|\hat{I}_n - I| \geq \varepsilon) \lesssim 2 - 2\Phi\left(\frac{\sqrt{n}\varepsilon}{\sqrt{\pi f^2 - (\pi f)^2}}\right),$$

and effectively computed using a consistent estimate or an upper bound of  $\pi f^2 - (\pi f)^2$ .

Assume now the MCMC setting, where no efficient procedure for sampling independent random variables from  $\pi$  is available. Let  $(X_n)_{n \geq 0}$  be an ergodic Markov chain on  $\mathcal{X}$  with transition kernel  $P$  and stationary limiting distribution  $\pi$ . Let  $\pi_0$  denote the initial distribution of the chain, i.e.  $X_0 \sim \pi_0$ . The distribution of  $X_t$  is  $\pi_t = \pi_0 P^t \rightarrow \pi$ , but  $X_0, X_1, \dots$  are dependent random variables and (1.2) is no longer an obvious and easy to analyze estimator. There are several possible strategies (cf. [Geyer 1992], [Niemi & Pokarowski 2007], [Chan & Yue 1996], [Liu, JS 2001], [Casella & Robert 1999]).

- *Estimation Along one Walk.* Use average along a single trajectory of the underlying Markov chain and discard the initial part to reduce bias. In this case the estimate is of the form

$$\hat{I}_{t,n} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i) \tag{1.3}$$

and  $t$  is called the burn-in time.

- *Estimation Along one Walk with Spacing.* Discard the initial part of a single trajectory to reduce bias and then take every  $s$ -th observation to reduce correlation. In this case the estimate is of the form

$$\hat{I}_{t,n,s} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_{is}) \quad (1.4)$$

and  $s$  is called the spacing parameter.

- *Multiple Run.* Use average over final states of multiple independent runs of the chain. Thus we need first to simulate say  $n$  trajectories of length say  $t$ :

$$\begin{array}{c} X_0^{(1)}, X_1^{(1)} \quad \dots, \quad X_t^{(1)}, \\ \vdots \\ X_0^{(n)}, X_1^{(n)} \quad \dots, \quad X_t^{(n)}, \end{array}$$

and for an estimate we take

$$\hat{I}_{t,n} = \frac{1}{n} \sum_{m=1}^n f(X_t^{(m)}), \quad (1.5)$$

where  $m$  numbers the independent runs of the chain and  $t$  should be large enough to reduce bias.

- *Median of Averages.* Use median of multiple independent shorter runs. Here we simulate
  - Simulate  $m$  independent runs of length  $t + n$  of the underlying Markov chain,

$$X_0^{(k)}, \dots, X_{t+n-1}^{(k)}, \quad k = 1, \dots, m.$$

- Calculate  $m$  estimates of  $I$ , each based on a single run,

$$\hat{I}_k = \hat{I}_{t,n}^{(k)} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i^{(k)}), \quad k = 1, \dots, m.$$

- For the final estimate take

$$\hat{I} = \text{med}\{\hat{I}_1, \dots, \hat{I}_m\}.$$

The one walk estimators are harder to analyze since both  $X_t, \dots, X_{t+n-1}$  and  $X_{ts}, \dots, X_{(t+n-1)s}$  are not independent, whereas  $X_t^{(1)}, \dots, X_t^{(m)}$  are. Yet one walk strategies are believed to be more efficient and are usually the practitioners' choice. Some precise results comparing the first three estimators under certain assumptions are available and confirm the practitioners' intuition. We refer to them later.

For each choice of estimation strategy additional questions arise, since one has to decide how to choose parameters  $t, n$  or  $t, n, s$  or  $t, n, m$  respectively, that assure "good quality of estimation". This choice must clearly depend on how one defines the desired "quality of estimation".

Moreover, we see from the above that MCMC requires a Markov chain on  $\mathcal{X}$  which is easily run on a computer, and which has  $\pi$  as its stationary limiting distribution. It may be a bit surprising that there exist reasonably general recipes for constructing such a chain that converges to  $\pi$  in most settings of practical interest.

### 1.2.1 The Metropolis Algorithm

The Metropolis algorithm has been introduced by Metropolis et al. in [Metropolis et al. 1953]. Let  $Q$  be a transition kernel of any other Markov chain that is easily simulated on a computer. Recall that  $\pi(\cdot)$  has a density  $\pi(dx) = p(x)dx$ , with possibly unknown normalizing constant. Let also  $Q(x, \cdot)$  have a density  $Q(x, dy) = q(x, y)dy$ . These densities are taken with respect to some  $\sigma$ -finite reference measure  $dx$ , which typically is the Lebesgue measure on  $R^d$ , however other settings are possible, including counting measures on discrete state spaces.

The Metropolis algorithm proceeds as follows.

1. Draw  $X_0$  from an initial distribution  $\pi_0$  (typically  $\pi_0 = \delta_{x_0}$  for some  $x_0 \in \mathcal{X}$ ).
2. Given  $X_n$  draw a proposal  $Y_{n+1}$  from  $Q(X_n, \cdot)$ .
3. Set

$$X_{n+1} = \begin{cases} Y_{n+1} & \text{with probability } \alpha(X_n, Y_{n+1}), \\ X_n & \text{with probability } 1 - \alpha(X_n, Y_{n+1}), \end{cases}$$

where

$$\alpha(x, y) := \min \left\{ 1, \frac{p(y)q(y, x)}{p(x)q(x, y)} \right\}$$

(Also, set  $\alpha(x, y) = 1$  whenever  $p(x)q(x, y) = 0$ .)

4. Replace  $n$  by  $n + 1$  and go to 2.

Note that one only has to compute the ratio of densities  $p(y)/p(x)$ , and hence the unknown normalizing constant for  $\pi$  in the acceptance probability  $\alpha(x, y)$  simplifies and one does not need to know it to run the chain.

Choosing the proposal density is another question that arises when implementing the Metropolis algorithm and different ways of doing it lead to different classes of algorithms. Typical classes include (see e.g. [Roberts & Rosenthal 2005])

- *Symmetric Metropolis Algorithm.* In this case  $q(x, y) = q(y, x)$  and hence  $\alpha(x, y) = \min\{1, \frac{\pi(y)}{\pi(x)}\}$ .
- *Random Walk Metropolis-Hastings.* In this case  $q(x, y) = q(y - x)$ .
- *Independence Sampler.* In this case the proposal does not depend on  $x$ , i.e.  $q(x, y) = q(y)$ .
- *Langevin Algorithm.* Where  $Q(X_n, \cdot) = N(X_n + (\delta/2)\nabla \log \pi(X_n), \delta)$  for some  $\delta > 0$ .

## 1.2.2 The Gibbs Sampler

The Gibbs Sampler is suitable in a setting where  $\mathcal{X}$  is a product space. For simplicity we suppose in this section that  $\mathcal{X}$  is an open subset of  $R^d$ , and write  $x = (x_1, \dots, x_d)$ .

The  $i$ -th component  $P_i$  of the Gibbs sampler  $P$  replaces  $x_i$  by a draw from the conditional distribution  $\pi(x_i|x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d)$ .

To state it more formally let, similarly as in [Roberts & Rosenthal 2005],

$$S_{x,i,a,b} = \{y \in \mathcal{X}; y_j = x_j \text{ for } j \neq i, \text{ and } a \leq y_i \leq b\}.$$

And

$$P_i(x, S_{x,i,a,b}) = \frac{\int_a^b p(x_1, \dots, x_{i-1}, t, x_{i+1}, \dots, x_d) dt}{\int_{-\infty}^{\infty} p(x_1, \dots, x_{i-1}, t, x_{i+1}, \dots, x_d) dt}. \quad (1.6)$$

Now the *deterministic scan Gibbs sampler* uses the transition kernel

$$P = P_1 P_2 \cdots P_d, \quad (1.7)$$

i.e. updates the coordinates of  $X_n$  in a systematic way, one after another, with draws from full conditional distributions.

On the other hand the *random scan Gibbs sampler* chooses a coordinate uniformly at random and performs its update, i.e. it uses the transition kernel

$$P = \frac{1}{d} \sum_{i=1}^d P_i. \quad (1.8)$$

In the example of Section 5.6 drawing from conditional distributions will be straightforward and in fact this is often the case for bayesian posterior distributions. However, if this step is infeasible, then instead of using  $P_i$  as defined in (1.6), one performs one step of a Metropolis algorithm designed to update  $i$ -th coordinate. Such a procedure is then called *Metropolis within Gibbs algorithm*.

### 1.3 Overview of the Results

Existing literature on Markov chains and their applications to Markov chain Monte Carlo procedures is to large extent focused on obtaining bounds on convergence rates to the stationary distribution ([Baxendale 2005], [Douc et al. 2003], [Jones & Hobert 2004], [Roberts & Tweedie 1999], [Rosenthal 1995b]) and on asymptotical results for MCMC estimators ([Jones et al. 2006], [Kipnis & Varadhan 1986], [Meyn & Tweedie 1993]). However, when analyzing MCMC estimators, results on the rate of convergence to the stationary distribution allow only to keep bias in control and do not translate in a straightforward way into bounds on the mean square error or confidence intervals. Moreover, asymptotic results may turn out useless in practice and may even be misleading ([Roberts & Rosenthal 2005]).

The main goal of this thesis is to obtain fixed-width results for an estimator, say  $\hat{I}$ , based on an MCMC algorithm. In particular we strive for the  $(\varepsilon - \alpha)$ -approximation, i.e.

$$P(|\hat{I} - I| \geq \varepsilon) \leq \alpha, \quad (1.9)$$

where  $\varepsilon$  is the desired quality of estimation and  $\alpha$  is the confidence level.

In analyzing Markov chains and estimators based on MCMC procedures we take the regenerative approach based on the split chain. The split chain construction allows to divide the Markov chain trajectory into independent

or 1–dependent blocks and turns out to be an extremely powerful technique with wide range of applications. The approach has been introduced independently in [Athreya & Ney 1978] and [Nummelin 1978] and immensely developed in [Nummelin 1984] and [Meyn & Tweedie 1993]. We give the basics of the approach in Chapter 2.

Results related to (1.9) are known in literature for discrete state space  $\mathcal{X}$  and bounded function  $f$  ([Aldous 1987], [Gillman 1998], [León & Perron 2004]). For general state space  $\mathcal{X}$ , and uniformly ergodic Markov chains (which in practice implies that  $\mathcal{X}$  is compact) and bounded function  $f$ , exponential inequalities are available (due to [Glynn & Ormoneit 2002] and an improved result due to [Kontoyiannis et al. 2005]) thus  $(\varepsilon - \alpha)$ –approximation can be easily deduced.

For a general, not necessarily compact, state space  $\mathcal{X}$  (or equivalently, not uniformly ergodic chains) and unbounded function  $f$  (which is e.g. the case when computing bayesian estimators for a quadratic loss function) no nonasymptotic results of type (1.9) are available. Fixed-width estimation is performed by deriving asymptotic confidence intervals based on

$$\hat{I}_n = \frac{1}{n} \sum_{i=0}^{n-1} f(X_i).$$

This construction requires two steps. First requirement is that a central limit theorem must hold, i.e.

$$\frac{\hat{I}_n - I}{\sqrt{n}} \xrightarrow{d} N(0, \sigma_f^2), \tag{1.10}$$

where  $\sigma_f^2 < \infty$  is the asymptotic variance. The second step is to obtain a strongly consistent estimator  $\hat{\sigma}_f^2$  of  $\sigma_f^2$ . Recent paper [Jones et al. 2006] presents the state of the art approach to the problem.

Results of Chapter 3 and Chapter 4 are related to this methodology.

In Chapter 3, based on [Bednorz, Latała & Łatuszyński 2008], a necessary and sufficient condition in terms of regeneration for a central limit theorem for functionals of ergodic Markov chains (as defined in (1.10)) have been obtained. It turns out, that the CLT holds if and only if excursions between regenerations are square integrable. An additional result of Chapter 3 is a solution to the open problem posed in [Roberts & Rosenthal 2005],

i.e. a regeneration proof of a CLT for uniformly ergodic Markov chains with  $E_\pi f^2 < \infty$ .

Chapter 4, based on [Bednorz & Łatuszyński 2007], is devoted to relaxing assumptions for strongly consistent estimators of  $\sigma_f^2$ . Results of Chapter 4 improve the methodology of [Jones et al. 2006].

In Chapter 5 nonasymptotic results of type (1.9) are obtained for noncompact state space  $\mathcal{X}$  and without assuming boundedness of the target function  $f$ .

More precisely, the goal of this chapter is to analyze estimation along one walk

$$\hat{I}_{t,n} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i) \quad (1.11)$$

of the unknown value  $I$  under the following drift condition towards a small set.

(A.1) Small set. There exist  $C \in \mathcal{B}(\mathcal{X})$ ,  $\tilde{\beta} > 0$  and a probability measure  $\nu$  on  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  such that for all  $x \in C$  and  $A \in \mathcal{B}(\mathcal{X})$

$$P(x, A) \geq \tilde{\beta}\nu(A).$$

(A.2) Drift. There exist a function  $V : \mathcal{X} \rightarrow [1, \infty)$  and constants  $\lambda < 1$  and  $K < \infty$  satisfying

$$PV(x) \leq \begin{cases} \lambda V(x), & \text{if } x \notin C, \\ K, & \text{if } x \in C. \end{cases}$$

(A.3) Aperiodicity. There exists  $\beta > 0$  such that  $\tilde{\beta}\nu(C) \geq \beta$ .

Under this assumption we provide explicit lower bounds on the *burn-in* time  $t$  and the length of simulation  $n$  that guarantee  $(\varepsilon - \alpha)$ -approximation. These bounds depend only and explicitly on the estimation parameters  $\varepsilon$  and  $\alpha$ , drift parameters  $\tilde{\beta}, \beta, \lambda, K$  and the the  $V$ -norm of the target function  $f$ , i.e.  $\|f^2\|_V = \sup_x f^2(x)/V(x)$ .

Moreover we analyze also estimation by the *median of averages* introduced in the previous section. It turns out that for small  $\alpha$  sharper bounds on the total simulation cost needed for  $(\varepsilon - \alpha)$ -approximation are available in this case by a simple exponential inequality.

The results of Chapter 5 have been applied for Gibbs samplers for a Hierarchical Random Effects Model of practical interest enabling nonasymptotic fixed-width analysis of this model. In particular this extends the results from [Jones & Hobert 2004], where burn in bounds in terms of total variation norm have been established for this model.

Chapter 6 deals with a slightly different topic, namely adaptive procedures. The idea is to modify the transition kernel based on the information collected during the simulation. This usually leads to a stochastic process that are not Markov chains any more and are less tractable theoretically. On the other hand, an adaptive procedure at time  $n$  is allowed to make use of an additional information: the sample trajectory up to time  $n$ . Clearly the class of stochastic processes used for simulation is bigger. Thus a smart use of the idea may lead to improvements in estimation quality. Simulations confirm this expectations and numerical examples for numerous specific algorithms outperform classical procedures [Roberts & Rosenthal 2006], [Kohn & Nott 2005]. An important example of the application of adaptive schemes is the Metropolis algorithm with multivariate normal proposal. In this case adaptation allows for automated choice of the covariance matrix for the proposal distribution [Atchadé & Rosenthal 2005]. Theoretical results on convergence and quality of estimation for adaptive procedures are very modest so far. Typical conditions that allow for investigation of convergence are called *diminishing adaptation* will be provided in Chapter 6. Time stability conditions for transition kernels assumed in ([Atchadé & Rosenthal 2005], [Kohn & Nott 2005]) fit into the diminishing adaptation framework. Intuitively time stability means that the adaptive process approaches a time homogeneous Markov chain.

In Chapter 6 we prove two results a convergence rate theorem and a law of large numbers for adaptive schemes. For both results we assume a *path stability condition* for transition kernels which is weaker than the *time stability condition*, assumed in [Atchadé & Rosenthal 2005] to prove similar results. The *path stability condition* results from *time stability condition* by the triangle inequality and intuitively means that the adaptive process approaches a time in-homogeneous Markov chain.



# Chapter 2

## Some Markov Chains

In this chapter we give some basic definitions and facts about stationarity and ergodicity of Markov chains that justify the Metropolis algorithm and the Gibbs sampler of Section 1.2 and provide grounds for the MCMC methodology. Next we outline the regeneration construction and the split chain and introduce typical objects and tools useful in for analyzing regenerative chains. Systematic, applications driven development of Markov chains theory via regeneration can be found e.g. in [Meyn & Tweedie 1993] and [Nummelin 1984] that constitute an immense body of work. Hence we do not attempt a systematic treatment of the Markov chain theory here and this chapter, based on [Meyn & Tweedie 1993], [Nummelin 1984], [Roberts & Rosenthal 2005] and [Nummelin 2002] is nothing more than a place for notions and tools frequently used in later chapters.

### 2.1 Stationarity and Ergodicity

Although majority of the results we describe carry over to the setting where  $\mathcal{X}$  is a general set and  $\mathcal{B}(\mathcal{X})$  is a countably generated  $\sigma$ -algebra (see e.g. [Meyn & Tweedie 1993]), in our applications driven development we believe Polish spaces offer more than sufficient generality and a great deal of "comfort". Thus, if not stated otherwise, the state space  $\mathcal{X}$  shall be a Polish and  $\mathcal{B}(\mathcal{X})$  shall denote the Borel  $\sigma$ -algebra on  $\mathcal{X}$ . A transition kernel  $P$  on  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  is a map  $P : \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$ , such that

- for any fixed  $A \in \mathcal{B}(\mathcal{X})$  the function  $P(\cdot, A)$  is measurable,

- for any fixed  $x \in \mathcal{X}$  the function  $P(x, \cdot)$  is a probability measure on  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ .

For a probability measure  $\mu$  and a transition kernel  $Q$ , by  $\mu Q$  we denote a probability measure defined by

$$\mu Q(\cdot) := \int_{\mathcal{X}} Q(x, \cdot) \mu(dx),$$

furthermore if  $g$  is a real-valued measurable function on  $\mathcal{X}$  let

$$Qg(x) := \int_{\mathcal{X}} g(y) Q(x, dy)$$

and

$$\mu g := \int_{\mathcal{X}} g(x) \mu(dx).$$

We will also use  $E_{\mu}g$  for  $\mu g$ , especially if  $\mu = \delta_x$  we will write  $E_x g$ . For transition kernels  $Q_1$  and  $Q_2$ ,  $Q_1 Q_2$  is also a transition kernel defined by

$$Q_1 Q_2(x, \cdot) := \int_{\mathcal{X}} Q_2(y, \cdot) Q_1(x, dy).$$

Let  $(X_n)_{n \geq 0}$  denote a time homogeneous Markov chain on  $\mathcal{X}$  evolving according to the transition kernel  $P$ , i.e. such that  $\mathcal{L}(X_{n+1}|X_n) = P(X_n, \cdot)$ . By  $\pi_0$  denote the distribution of  $X_0$ , i.e. the initial distribution of the chain. Then, using the above notation the distribution of  $X_n$  is  $\pi_n = \pi_0 P^n$ . In particular, if  $\pi = \delta_x$ , then  $X_n$  is distributed as  $\pi_n = \delta_x P^n = P^n(x, \cdot)$ . Clearly the behavior of  $\pi_n$  is of our vital interest.

We say that a probability distribution  $\pi$  is stationary for  $P$ , if  $\pi P = \pi$ . A crucial notion related to stationarity via Proposition 2.1.2 is reversibility.

**Definition 2.1.1.** A Markov chain on a state space  $\mathcal{X}$  with transition kernel  $P$  is reversible with respect to a probability distribution  $\pi$  on  $\mathcal{X}$ , if

$$\int_A P(x, B) \pi(dx) = \int_B P(y, A) \pi(dy), \quad \text{for all } A, B \in \mathcal{B}(\mathcal{X})$$

we shall write equivalently

$$\pi(dx) P(x, dy) = \pi(dy) P(y, dx), \quad \text{for all } x, y \in \mathcal{X}.$$

**Proposition 2.1.2.** *If a Markov chain with transition kernel  $P$  is reversible with respect to  $\pi$ , then  $\pi$  is stationary for  $P$ .*

*Proof.*

$$\pi P(A) = \int_{\mathcal{X}} P(x, A) \pi(dx) = \int_A P(y, \mathcal{X}) \pi(dy) = \int_A \pi(dy) = \pi(A).$$

□

It is straightforward to check that the acceptance probability  $\alpha(x, y)$  of the Metropolis algorithm of Section 1.2.1 makes the procedure reversible with respect to  $\pi$  and thus it has  $\pi$  as its stationary distribution.

Also the  $i$ -th component  $P_i$  of the Gibbs sampler of Section 1.2.2 is a special case of the Metropolis algorithm (with  $\alpha(x, y) = 1$ ) and hence  $\pi$  is stationary for  $P_i$ . This implies that the random scan Gibbs sampler is reversible and has  $\pi$  as its stationary distribution. The deterministic scan Gibbs sampler usually is not reversible, however since  $\pi$  is stationary for each  $P_i$ , it is also stationary for  $P$ .

Obviously stationarity is not enough for the applications in question since it does not even imply  $\pi_n \rightarrow \pi$  (see [Roberts & Rosenthal 2005] for examples), not to mention justifying any of the estimation schemes (1.3-1.5). One needs some more assumptions and notions to investigate convergence of  $\pi_n$  to  $\pi$  and properties of estimation strategies of previous sections.

In particular the total variation distance is a very common tool to evaluate distance between two probability measures  $\mu_1$  and  $\mu_2$  and is defined as

$$\|\mu_1 - \mu_2\|_{tv} = \sup_{A \in \mathcal{B}(\mathcal{X})} |\mu_1(A) - \mu_2(A)|. \quad (2.1)$$

We shall distinguish between the two following types of convergence to  $\pi$ .

$$\lim_{n \rightarrow \infty} \|P^n(x, \cdot) - \pi\|_{tv} = 0, \quad \text{for } \pi\text{-almost every } x \in \mathcal{X}, \quad (2.2)$$

$$\lim_{n \rightarrow \infty} \|P^n(x, \cdot) - \pi\|_{tv} = 0, \quad \text{for all } x \in \mathcal{X}. \quad (2.3)$$

$\phi$ -irreducibility and aperiodicity are properties that guarantee convergence in (2.2).

**Definition 2.1.3.** A Markov chain  $(X)_{n \geq 0}$  with transition kernel  $P$  is  $\phi$ -irreducible if there exists a non-zero  $\sigma$ -finite measure  $\phi$  on  $\mathcal{X}$  such that for all  $A \subseteq \mathcal{X}$  with  $\phi(A) > 0$ , and for all  $x \in \mathcal{X}$ , there exists a positive integer  $n = n(x, A)$  such that  $P^n(x, A) > 0$ .

**Definition 2.1.4.** A Markov chain  $(X)_{n \geq 0}$  with transition kernel  $P$  and stationary distribution  $\pi$  is periodic with period  $d \geq 2$  if there exist disjoint subsets  $\mathcal{X}_0, \dots, \mathcal{X}_{d-1} \subseteq \mathcal{X}$  such that  $\pi(\mathcal{X}_1) > 0$  and for all  $0 \leq i \leq d-1$ , and for all  $x \in \mathcal{X}_i$ ,  $P(x, \mathcal{X}_{i+1 \bmod d}) = 1$ . And  $d$  is maximal for the property. Otherwise the chain is called aperiodic.

**Theorem 2.1.5.** *If a Markov chain  $(X)_{n \geq 0}$  with transition kernel  $P$  and stationary distribution  $\pi$  on a state space  $\mathcal{X}$  is  $\phi$ -irreducible and aperiodic, then (2.2) holds.*

*Moreover, if a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is such that  $\pi(|f|) < \infty$ , then a strong law of large numbers holds in the following sense*

$$\frac{1}{n} \sum_{i=0}^{n-1} f(X_i) \rightarrow \pi f, \quad \text{as } n \rightarrow \infty, \quad \text{w.p. 1.} \quad (2.4)$$

The foregoing convergence result is one of many possible formulations. A proof of the first part can be found in [Roberts & Rosenthal 2005] Section 4.6 and the strong law of large numbers part results e.g. from Theorem 17.0.1 of [Meyn & Tweedie 1993]. Theorem 2.1.5 is widely applicable to MCMC algorithms. The Metropolis algorithm and the Gibbs samplers of Section 1.2 are designed precisely so that  $\pi$  is stationary. Also, it is usually straightforward to verify that the chain is aperiodic and  $\phi$ -irreducible with e.g.  $\phi$  being the Lebesgue measure or  $\phi = \pi$ .

The following example due to C. Geyer (cf. [Roberts & Rosenthal 2005]) provides a simple Markov chain that exhibits a "bad" behavior on a null set.

**Example 2.1.6.** Let  $\mathcal{X} = \{1, 2, \dots\}$  and define transition probabilities by  $P(1, \{1\}) = 1$ , and for  $x \geq 2$ , let  $P(x, \{1\}) = 1/x^2$  and  $P(x, \{x+1\}) = 1 - 1/x^2$ . Then the chain is aperiodic and  $\pi = \delta_1$  is the invariant distribution. The chain is also  $\pi$ -irreducible. However, if  $X_0 = x \geq 2$ , then  $P(X_n = x+n \text{ for all } n) > 0$ , and  $\|P^n(x, \cdot) - \pi(\cdot)\| \not\rightarrow 0$ . Thus the convergence holds only for  $x = 1$  which in this case is  $\pi$ -a.e.  $x \in \mathcal{X}$ .

To guarantee convergence for all  $x \in \mathcal{X}$ , as in (2.3) one needs to assume slightly more, namely Harris recurrence.

**Definition 2.1.7** (Harris Recurrence). A Markov chain  $(X_n)_{n \geq 0}$  with transition kernel  $P$  and stationary probability measure  $\pi$  is Harris recurrent if for all  $A \in \mathcal{B}(\mathcal{X})$ , such that  $\pi(A) > 0$ , and all  $x \in \mathcal{X}$ , the chain started at  $x$  will eventually reach  $A$  with probability 1, i.e.  $P(\exists n : X_n \in A | X_0 = x) = 1$ .

**Theorem 2.1.8.** *Ergodicity as defined in (2.3) is equivalent to Harris recurrence and aperiodicity.*

The foregoing Theorem 2.1.8 results from Proposition 6.3 in [Nummelin 1984]. Harris recurrent and aperiodic chains are often referred to as *Harris ergodic*.

The speed of convergence in (2.2) or (2.3) is another natural criterion for classifying chains. Geometrically ergodic and uniformly ergodic chains are of particular interest.

**Definition 2.1.9** (Uniform Ergodicity and Geometric Ergodicity). We say that a Markov chain  $(X_n)_{n \geq 0}$  with transition kernel  $P$  and stationary distribution  $\pi$  is

- *geometrically ergodic*, if  $\|P^n(x, \cdot) - \pi(\cdot)\|_{tv} \leq M(x)\rho^n$ , for some  $\rho < 1$  and  $M(x) < \infty$   $\pi$ -almost everywhere,
- *uniformly ergodic*, if  $\|P^n(x, \cdot) - \pi(\cdot)\|_{tv} \leq M\rho^n$ , for some  $\rho < 1$  and  $M < \infty$ ,

The difference between geometric ergodicity and uniform ergodicity is that  $M$  may depend on the initial state  $x$ . Obviously, if a chain is geometrically ergodic and  $M(x)$  is a bounded function, then the chain is also uniformly ergodic. In particular, if the state space is finite, then every geometrically ergodic Markov chain is uniformly ergodic. (And from the standard theory of discrete state space Markov chains we know that every ergodic chain is uniformly ergodic.) Verifying uniform or geometric ergodicity is in general nontrivial and we will refer to it later. An interesting result for the algorithms presented in Chapter 1 is for example that a symmetric random-walk Metropolis algorithm is geometrically ergodic if and only if  $\pi$  has finite exponential moments, as shown in [Mengersen & Tweedie 1996].

Since in the sequel we deal with integrals of unbounded functions  $f$  with respect to probability measures, the very common total variation distance defined by (2.1) is in this case inappropriate for measuring distances between probability measures and we need to introduce the  $V$ -norm and  $V$ -norm distance.

Let  $V : \mathcal{X} \rightarrow [1, \infty)$  be a measurable function. For measurable function  $g : \mathcal{X} \rightarrow \mathbb{R}$  define its  $V$ -norm as

$$|g|_V := \sup_{x \in \mathcal{X}} \frac{|g(x)|}{V(x)}.$$

To evaluate the distance between two probability measures  $\mu_1$  and  $\mu_2$  we use the *V-norm distance*, defined for probability measures  $\mu_1$  and  $\mu_2$  as

$$\|\mu_1 - \mu_2\|_V := \sup_{|g| \leq V} |\mu_1 g - \mu_2 g|.$$

Note that for  $V \equiv 1$  the  $V$ -norm distance  $\|\cdot\|_V$  amounts to the total variation distance, i.e.  $\|\mu_1 - \mu_2\|_V = 2 \sup_{A \in \mathcal{B}(\mathcal{X})} |\mu_1(A) - \mu_2(A)| = 2\|\mu_1 - \mu_2\|_{tv}$ . Finally for two transition kernels  $Q_1$  and  $Q_2$  the *V-norm distance* between  $Q_1$  and  $Q_2$  is defined by

$$\|Q_1 - Q_2\|_V := \|Q_1(x, \cdot) - Q_2(x, \cdot)\|_V|_V = \sup_{x \in \mathcal{X}} \frac{\|Q_1(x, \cdot) - Q_2(x, \cdot)\|_V}{V(x)}.$$

For a probability distribution  $\mu$ , define a transition kernel  $\mu(x, \cdot) := \mu(\cdot)$ , to allow for writing  $\|Q - \mu\|_V$  and  $\|\mu_1 - \mu_2\|_V$ . Define also the following Banach space

$$B_V := \{f : f : \mathcal{X} \rightarrow R, |f|_V < \infty\}.$$

Now if  $\|Q_1 - Q_2\|_V < \infty$ , then  $Q_1 - Q_2$  is a bounded operator from  $B_V$  to itself, and  $\|Q_1 - Q_2\|_V$  is its operator norm. See [Meyn & Tweedie 1993] Chapter 16 for details.

Now we are in a position to introduce the  $V$ -uniform ergodicity.

**Definition 2.1.10** (*V*-uniform ergodicity). We say that a Markov chain  $(X_n)_{n \geq 0}$  with transition kernel  $P$  and stationary distribution  $\pi$  is *V*-uniformly ergodic, if

$$\|P^n - \pi\|_V \rightarrow 0, \quad \text{as } n \rightarrow \infty. \quad (2.5)$$

Moreover, since  $\|\cdot\|_V$  is an operator norm (2.5) is equivalent to

$$\|P^n - \pi\|_V \leq M\rho^n, \quad \text{for some } M < \infty \text{ and } \rho < 1. \quad (2.6)$$

## 2.2 Small Sets and the Split Chain

The regeneration construction has been invented independently by [Nummelin 1978] and [Athreya & Ney 1978] and is now a very celebrated technique. The development of this approach resulted in intuitive and rather simple proofs of most results about Markov chains and enabled better understanding and rapid progress of the theory. In this section we provide the basics of the

regeneration and split chain construction needed for the following chapters. Systematic development of the theory can be found in [Nummelin 1984] and [Meyn & Tweedie 1993] which we exploit here.

We begin with the following definition of an atom.

**Definition 2.2.1** (Atom). A set  $B \in \mathcal{B}(\mathcal{X})$  is called an atom for a Markov chain  $(X)_{n \geq 0}$  with transition kernel  $P$  if there exists a probability measure  $\nu$  on  $\mathcal{B}(\mathcal{X})$ , such that for all  $x \in B$ ,

$$P(x, \cdot) = \nu(\cdot).$$

If the Markov chain is  $\psi$ -irreducible and  $\psi(B) > 0$  then  $B$  is called an accessible atom.

A single point  $x \in \mathcal{X}$  is always an atom. For a discrete state space irreducible Markov chain every single point is an accessible atom. Much of the discrete state space theory is developed by studying Markov chain tours between consecutive visits to a distinguished atom  $c \in \mathcal{X}$ . On a general state space accessible atoms typically do not exist. However such atoms can be artificially constructed. First we provide a general version of a minorization condition that enables this construction.

**Definition 2.2.2** (Minorization Condition - a general version). Let  $s : \mathcal{X} \rightarrow [0, 1]$  be a function for which  $E_\pi s > 0$  and there exists an  $m > 0$  and such a probability measure  $\nu_m$  on  $\mathcal{B}(\mathcal{X})$ , that for all  $x \in \mathcal{X}$ ,

$$P^m(x, \cdot) \geq s(x)\nu_m(\cdot). \quad (2.7)$$

However, a special case of this condition with  $s(x) = \varepsilon \mathbb{1}_C(x)$  usually turns out to be as powerful as the general version and is often more suitable to work with.

**Definition 2.2.3** (Small Set). A set  $C \in \mathcal{B}(\mathcal{X})$  is  $\nu_m$ -small, if there exist  $m > 0$ ,  $\varepsilon > 0$ , and a probability measure  $\nu_m$  on  $\mathcal{B}(\mathcal{X})$ , such that for all  $x \in C$ ,

$$P^m(x, \cdot) \geq \varepsilon \nu_m(\cdot). \quad (2.8)$$

*Remark 2.2.4.* Theorem 5.2.2 of [Meyn & Tweedie 1993] states that any  $\psi$ -irreducible Markov chain is well-endowed with small sets  $C$  of positive measure  $\psi$  and such that  $\nu_m(C) > 0$ . Since ergodic Markov chains are  $\pi$ -irreducible, for an ergodic chain a small set  $C$  with  $\pi(C) > 0$  and  $\nu_m(C) > 0$  always exists.

Definition 2.2.3 and Remark 2.2.4 imply the following minorization condition.

**Definition 2.2.5** (Minorization Condition). For some  $\varepsilon > 0$ , some  $C$  such that  $\psi(C) > 0$ , and some probability measure  $\nu_m$  with  $\nu_m(C) = 1$  we have for all  $x \in C$ ,

$$P^m(x, \cdot) \geq \varepsilon \nu_m(\cdot). \quad (2.9)$$

The minorization condition (2.9) allows for constructing the split chain for  $(X_n)_{n \geq 0}$  which is the central object of the approach (see Section 17.3 of [Meyn & Tweedie 1993] for a detailed description). Let  $(X_{nm})_{n \geq 0}$  be the  $m$ -skeleton of  $(X_n)_{n \geq 0}$ , i.e. a Markov chain evolving according to the  $m$ -step transition kernel  $P^m$ . The minorization condition allows to write  $P^m$  as a mixture of two distributions:

$$P^m(x, \cdot) = \varepsilon \mathbb{I}_C(x) \nu_m(\cdot) + [1 - \varepsilon \mathbb{I}_C(x)] R(x, \cdot), \quad (2.10)$$

where  $R(x, \cdot) = [1 - \varepsilon \mathbb{I}_C(x)]^{-1} [P(x, \cdot) - \varepsilon \mathbb{I}_C(x) \nu_m(\cdot)]$ . Now let  $(X_{nm}, Y_n)_{n \geq 0}$  be the split chain of the  $m$ -skeleton i.e. let the random variable  $Y_n \in \{0, 1\}$  be the level of the split  $m$ -skeleton at time  $nm$ . The split chain  $(X_{nm}, Y_n)_{n \geq 0}$  is a Markov chain that obeys the following transition rule  $\check{P}$ .

$$\check{P}(Y_n = 1, X_{(n+1)m} \in dy | Y_{n-1}, X_{nm} = x) = \varepsilon \mathbb{I}_C(x) \nu_m(dy) \quad (2.11)$$

$$\check{P}(Y_n = 0, X_{(n+1)m} \in dy | Y_{n-1}, X_{nm} = x) = (1 - \varepsilon \mathbb{I}_C(x)) R(x, dy), \quad (2.12)$$

and  $Y_n$  can be interpreted as a coin toss indicating whether  $X_{(n+1)m}$  given  $X_{nm} = x$  should be drawn from  $\nu_m(\cdot)$  - with probability  $\varepsilon \mathbb{I}_C(x)$  - or from  $R(x, \cdot)$  - with probability  $1 - \varepsilon \mathbb{I}_C(x)$ .

Obviously  $(X_{nm}, Y_n)_{n \geq 0}$ , i.e. the split chain of the  $m$ -skeleton is a Markov chain and the crucial observation follows from the Bayes rule, namely the set  $\check{\alpha} := C \times \{1\}$  is an accessible atom for this chain.

One obtains the split chain  $(X_k, Y_n)_{k \geq 0, n \geq 0}$  of the initial Markov chain  $(X_n)_{n \geq 0}$  by defining appropriate conditional probabilities. To this end let  $X_0^{nm} = \{X_0, \dots, X_{nm-1}\}$  and  $Y_0^n = \{Y_0, \dots, Y_{n-1}\}$ .

$$\check{P}(Y_n = 1, X_{nm+1} \in dx_1, \dots, X_{(n+1)m-1} \in dx_{m-1}, X_{(n+1)m} \in dy | \quad (2.13)$$

$$|Y_0^n, X_0^{nm}; X_{nm} = x) = \frac{\varepsilon \mathbb{I}_C(x) \nu_m(dy)}{P^m(x, dy)} P(x, dx_1) \cdots P(x_{m-1}, dy),$$

$$\check{P}(Y_n = 0, X_{nm+1} \in dx_1, \dots, X_{(n+1)m-1} \in dx_{m-1}, X_{(n+1)m} \in dy | \quad (2.14)$$

$$|Y_0^n, X_0^{nm}; X_{nm} = x) = \frac{(1 - \varepsilon \mathbb{I}_C(x)) R(x, dy)}{P^m(x, dy)} P(x, dx_1) \cdots P(x_{m-1}, dy),$$



where  $\frac{\nu_m(dy)}{P^m(x,dy)}$  and  $\frac{R(x,dy)}{P^m(x,dy)}$  are Radon-Nykodym derivatives. Note that the marginal distribution of  $(X_k)_{k \geq 0}$  in the split chain is that of the underlying Markov chain with transition kernel  $P$ .

An important characterization of the invariant measure obtained via the splitting technique is a generalization of the Kac's Theorem, namely Theorem 2.2.8, which is the key conclusion of Chapter 10 in [Meyn & Tweedie 1993]. Let

$$U(x, A) := \sum_{n=1}^{\infty} P^n(x, A) = E_x \left( \sum_{n=1}^{\infty} \mathbb{I}_A(X_n) \right)$$

and for a measure  $\psi$  define

$$\mathcal{B}^+(\mathcal{X}) := \{A \in \mathcal{B}(\mathcal{X}) : \psi(A) > 0\}.$$

**Definition 2.2.6** (Recurrent Chains). A chain  $(X_n)_{n \geq 0}$  with a transition kernel  $P$  is called recurrent if it is  $\psi$ -irreducible and  $U(x, A) = \infty$  for any  $x \in \mathcal{X}$  and every  $A \in \mathcal{B}^+(\mathcal{X})$ .

*Remark 2.2.7.* Recurrence is a weaker condition than Harris recurrence, in particular the Markov chain defined in Example 2.1.6 is recurrent but not Harris recurrent.

Moreover, for a set  $A \in \mathcal{X}$  define its hitting time  $\tau_A$  as

$$\tau_A := \min\{n \geq 1 : X_n \in A\}.$$

**Theorem 2.2.8.** *Let the Markov chain  $(X_n)_{n \geq 0}$  be recurrent. Then there exists a unique (up to constant multiples) invariant measure  $\pi_u$ . This measure  $\pi_u$  has the following representation for any  $A \in \mathcal{B}^+(\mathcal{X})$*

$$\pi_u(B) = \int_A E_x \left[ \sum_{n=1}^{\tau_A} \mathbb{I}_B(X_n) \right] \pi_u(dx), \quad B \in \mathcal{B}(\mathcal{X}). \quad (2.15)$$

Moreover, the measure  $\pi_u$  is finite if there exists a small set  $C$  such that

$$\sup_{x \in C} E_x[\tau_C] < \infty.$$

To take advantage of the splitting technique for analyzing Markov chains and functionals of Markov chains we need a bit more formalism. For a measure  $\lambda$  on  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  let  $\lambda^*$  denote the measure on  $\mathcal{X} \times \{0, 1\}$  (with

product  $\sigma$ -algebra) defined by  $\lambda^*(B \times \{1\}) = \varepsilon\lambda(B \cap C)$  and  $\lambda^*(B \times \{0\}) = (1 - \varepsilon)\lambda(B \cap C) + \lambda(B \cap C^c)$ . In the sequel we shall use  $\nu_m^*$  for which  $\nu_m^*(B \times \{1\}) = \varepsilon\nu_m(B)$  and  $\nu_m^*(B \times \{0\}) = (1 - \varepsilon)\nu_m(B)$  due to the fact that  $\nu_m(C) = 1$ .

Now integrate (2.13) over  $x_1, \dots, x_{m-1}$  and then over  $y$ . This yields

$$\check{P}(Y_n = 1, X_{(n+1)m} \in dy | Y_0^n, X_0^{nm}; X_{nm} = x) = \varepsilon \mathbb{I}_C(x) \nu_m(dy), \quad (2.16)$$

and

$$\check{P}(Y_n = 1 | Y_0^n, X_0^{nm}; X_{nm} = x) = \varepsilon \mathbb{I}_C(x). \quad (2.17)$$

From the Bayes rule we obtain

$$\check{P}(X_{(n+1)m} \in dy | Y_0^n, X_0^{nm}; Y_n = 1, X_{nm} = x) = \nu_m(dy), \quad (2.18)$$

and the crucial observation due to Meyn and Tweedie, emphasized here as Lemma 2.2.9 follows.

**Lemma 2.2.9.** *Conditional on  $\{Y_n = 1\}$ , the pre- $nm$  process  $\{X_k, Y_i : k \leq nm, i \leq n\}$  and the post- $(n+1)m$  process  $\{X_k, Y_i : k \geq (n+1)m, i \geq n+1\}$  are independent. Moreover, the post- $(n+1)m$  process has the same distribution as  $\{X_k, Y_i : k \geq 0, i \geq 0\}$  with  $\nu_m^*$  for the initial distribution of  $(X_0, Y_0)$ .*

Next, let  $\sigma_{\check{\alpha}}(n)$  denote entrance times of the split chain to the set  $\check{\alpha} = C \times \{1\}$ , i.e.

$$\sigma_{\check{\alpha}}(0) = \min\{k \geq 0 : Y_k = 1\}, \quad \sigma_{\check{\alpha}}(n) = \min\{k > \sigma_{\check{\alpha}}(n-1) : Y_k = 1\}, \quad n \geq 1,$$

whereas hitting times  $\tau_{\check{\alpha}}(n)$  are defined as follows:

$$\tau_{\check{\alpha}}(1) = \min\{k \geq 1 : Y_k = 1\}, \quad \tau_{\check{\alpha}}(n) = \min\{k > \tau_{\check{\alpha}}(n-1) : Y_k = 1\}, \quad n \geq 2.$$

In view of Lemma 2.2.9 it should be intuitively clear that the following tours

$$\left\{ \left\{ X_{(\sigma_{\check{\alpha}}(n)+1)m}, X_{(\sigma_{\check{\alpha}}(n)+1)m+1}, \dots, X_{(\sigma_{\check{\alpha}}(n+1)+1)m-1} \right\}, n = 0, 1, \dots \right\}$$

that start whenever  $X_k \sim \nu_m$  are of crucial importance. In fact in the next chapter they will turn out to be much more tractable than the crude chain  $(X_n)_{n \geq 0}$  on  $\mathcal{X}$ .

Since we are interested in functionals of the Markov chain  $(X_n)_{n \geq 0}$ , for a real-valued function, say  $g$ , on  $\mathcal{X}$ , we define here also

$$s_i = s_i(g) = \sum_{j=m(\sigma_{\bar{\alpha}}(i)+1)}^{m(\sigma_{\bar{\alpha}}(i+1)+1)-1} g(X_j) = \sum_{j=\sigma_{\bar{\alpha}}(i)+1}^{\sigma_{\bar{\alpha}}(i+1)} Z_j(g), \quad (2.19)$$

where

$$Z_j(g) = \sum_{k=0}^{m-1} g(X_{jm+k}). \quad (2.20)$$

*Remark 2.2.10.* Clearly, one can construct the split chain based on the more general minorization condition (2.7) instead of (2.9). We chose (2.9) for simplicity. However, we use the split chain construction based on (2.7) in Chapter 4.

# Chapter 3

## A Complete Characterisation of $\sqrt{n}$ -CLTs for Ergodic Markov Chains via Regeneration

Central limit theorems for functionals of general state space Markov chains are of crucial importance in sensible implementation of Markov chain Monte Carlo algorithms as well as of vital theoretical interest. Different approaches to proving this type of results under diverse assumptions led to a large variety of CLT versions. However due to the recent development of the regeneration theory of Markov chains, many classical CLTs can be reproved using this intuitive probabilistic approach, avoiding technicalities of original proofs. In this paper we provide an if and only if characterization of  $\sqrt{n}$ -CLTs for ergodic Markov chains via regeneration and then use the result to solve the open problem posed in [Roberts & Rosenthal 2005]. We then discuss the difference between one-step and multiple-step small set condition.

Results of this chapter are based on paper [Bednorz, Latała & Łatuszyński 2008] and are joint work with Witold Bednorz and Rafał Latała.

### 3.1 CLTs for Markov Chains

Let  $(X_n)_{n \geq 0}$  be a time homogeneous, ergodic Markov chain on a measurable space  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ , with transition kernel  $P$  and a unique stationary measure  $\pi$  on  $\mathcal{X}$ . We remark that here ergodicity means that

$$\lim_{n \rightarrow \infty} \|P^n(x, \cdot) - \pi\|_{tv} = 0, \quad \text{for all } x \in \mathcal{X}, \quad (3.1)$$

where  $\|\cdot\|_{tv}$  denotes the total variation distance. The process  $(X_n)_{n \geq 0}$  may start from any initial distribution  $\pi_0$ . Let  $g$  be a real valued Borel function on  $\mathcal{X}$ , square integrable against the stationary measure  $\pi$ . We denote by  $\bar{g}$  its centered version, namely  $\bar{g} = g - \int g d\pi$  and for simplicity  $S_n := \sum_{i=0}^{n-1} \bar{g}(X_i)$ . We say that a  $\sqrt{n}$ -CLT holds for  $(X_n)_{n \geq 0}$  and  $g$  if

$$S_n/\sqrt{n} \xrightarrow{d} N(0, \sigma_g^2), \quad \text{as } n \rightarrow \infty, \quad (3.2)$$

where  $\sigma_g^2 < \infty$ .

Central limit theorems as defined by condition (3.2) are crucial for assessing the quality of Markov chain Monte Carlo estimation as we demonstrate in Chapter 4 (c.f. [Jones et al. 2006] and [Geyer 1992]) and are also of independent theoretical interest. Thus a large body of work on CLTs for functionals of Markov chains exists and a variety of results have been established under different assumptions and with different approaches to proofs (see [Jones 2005] for a review).

First we aim to provide a general result, namely Theorem 3.3.1, that gives a necessary and sufficient condition for  $\sqrt{n}$ -CLTs for ergodic chains (which is a generalization of the well known Theorem 17.3.6 [Meyn & Tweedie 1993]). Assume for a moment that there exists an accessible atom  $\alpha \in \mathcal{B}(\mathcal{X})$ , i.e. such a set  $\alpha$  that  $\pi(\alpha) > 0$  and there exists a probability measure  $\nu$  on  $\mathcal{B}(\mathcal{X})$ , such that  $P(x, A) = \nu(A)$  for all  $x \in \alpha$ . Let  $\tau_\alpha$  be the first hitting time for  $\alpha$ . In this simplistic case we can rephrase our Theorem 3.3.1 as follows:

**Theorem 3.1.1.** *Suppose that  $(X_n)_{n \geq 0}$  is ergodic and possess an accessible atom  $\alpha$ , then the  $\sqrt{n}$ -CLT holds if and only if*

$$E_\alpha \left[ \left( \sum_{k=1}^{\tau_\alpha} \bar{g}(X_k) \right)^2 \right] < \infty. \quad (3.3)$$

Furthermore we have the following formula for the variance

$$\sigma_g^2 = \pi(\alpha) E_\alpha \left[ \left( \sum_{k=1}^{\tau_\alpha} \bar{g}(X_k) \right)^2 \right].$$

We discuss briefly the relation between two classical CLT formulations for geometrically ergodic and uniformly ergodic Markov chains (recall Definition 2.1.9). Recently the following CLT provided by [Ibragimov & Linnik 1971] has been reproved in [Roberts & Rosenthal 2005] using the intuitive regeneration approach and avoiding technicalities of the original proof (however see Section 3.5 for a commentary).

**Theorem 3.1.2.** *If a Markov chain  $(X_n)_{n \geq 0}$  with stationary distribution  $\pi$  is geometrically ergodic, then a  $\sqrt{n}$ -CLT holds for  $(X_n)_{n \geq 0}$  and  $g$  whenever  $\pi(|g|^{2+\delta}) < \infty$  for some  $\delta > 0$ . Moreover  $\sigma_g^2 := \int_{\mathcal{X}} \bar{g}^2 d\pi + 2 \int_{\mathcal{X}} \sum_{n=1}^{\infty} \bar{g}(X_0) \bar{g}(X_n) d\pi$ .*

*Remark 3.1.3.* Note that for reversible chains the condition  $\pi(|g|^{2+\delta}) < \infty$  for some  $\delta > 0$  in Theorem 3.1.2 can be weakened to  $\pi(g^2) < \infty$  as proved in [Roberts & Rosenthal 1997b], however this is not possible for the general case, see [Bradley 1983] or [Häggström 2005] for counterexamples.

Roberts and Rosenthal posed an open problem, whether the following CLT version for uniformly ergodic Markov chains due to [Cogburn 1972] can also be reproved using direct regeneration arguments.

**Theorem 3.1.4.** *If a Markov chain  $(X_n)_{n \geq 0}$  with stationary distribution  $\pi$  is uniformly ergodic, then a  $\sqrt{n}$ -CLT holds for  $(X_n)_{n \geq 0}$  and  $g$  whenever  $\pi(g^2) < \infty$ . Moreover  $\sigma_g^2 := \int_{\mathcal{X}} \bar{g}^2 d\pi + 2 \int_{\mathcal{X}} \sum_{n=1}^{\infty} \bar{g}(X_0) \bar{g}(X_n) d\pi$ .*

The aim of this chapter is to prove Theorem 3.3.1 and show how to derive from this general framework the regeneration proof of Theorem 3.1.4. The outline of the chapter is as follows. In Section 3.2 we provide some preliminary results which may also be of independent interest. In Section 3.3 we detail the proof of Theorem 3.3.1, and derive Theorem 3.1.4 as a corollary in Section 3.4. Section 3.5 comprises a discussion of some difficulties of the regeneration approach.

## 3.2 Tools and Preliminary Results

Recall the split chain construction of the previous chapter and the notation therein. In particular  $s_i$ , defined by (2.19) will be of our vital interest.

In this section we take  $\bar{g}$ , the centered version of  $g$ , and analyze the sequence  $s_i(\bar{g})$ ,  $i \geq 0$ . The basic result we often refer to is Theorem 17.3.1 in [Meyn & Tweedie 1993], which states that  $(s_i)_{i \geq 0}$  is a sequence of 1-dependent, identically distributed r.v.'s with  $\check{E}s_i = 0$ . In our approach we use the following decomposition:  $s_i = \underline{s}_i + \bar{s}_i$ , where

$$\underline{s}_i := \sum_{j=\sigma_{\bar{\alpha}}(i)+1}^{\sigma_{\bar{\alpha}}(i+1)-1} Z_j(\bar{g}) - \check{E}_{\pi_0^*} \left[ \sum_{j=\sigma_{\bar{\alpha}}(i)+1}^{\sigma_{\bar{\alpha}}(i+1)-1} Z_j(\bar{g}) \right], \quad (3.4)$$

$$\bar{s}_i := Z_{\sigma_{\bar{\alpha}}(i+1)}(\bar{g}) - \check{E}_{\pi_0^*} \left[ Z_{\sigma_{\bar{\alpha}}(i+1)}(\bar{g}) \right]. \quad (3.5)$$

A look into the proof of Lemma 3.2.3 later in this section clarifies that  $\underline{s}_i$  and  $\bar{s}_i$  are well defined.

**Lemma 3.2.1.** *The sequence  $(\underline{s}_i)_{i \geq 0}$  consists of i.i.d. random variables.*

*Proof.* First note that  $\underline{s}_i$  is a function of  $\{X_{(\sigma_{\bar{\alpha}}(i)+1)m}, X_{(\sigma_{\bar{\alpha}}(i)+1)m+1}, \dots\}$  and that  $Y_{\sigma_{\bar{\alpha}}(i)} = 1$ , hence by Lemma 2.2.9  $\underline{s}_0, \underline{s}_1, \underline{s}_2, \dots$  are identically distributed. Now focus on  $\underline{s}_i, \underline{s}_{i+k}$  and  $Y_{\sigma_{\bar{\alpha}}(i+k)}$  for some  $k \geq 1$ . Obviously  $Y_{\sigma_{\bar{\alpha}}(i+k)} = 1$ . Moreover  $\underline{s}_i$  is a function of the pre- $\sigma_{\bar{\alpha}}(i+k)m$  process and  $\underline{s}_{i+k}$  is a function of the post- $(\sigma_{\bar{\alpha}}(i+k)+1)m$  process. Thus  $\underline{s}_i$  and  $\underline{s}_{i+k}$  are independent again by Lemma 2.2.9 and for  $A_i, A_{i+k}$ , Borel subsets of  $R$ , we have

$$\check{P}_{\pi_0^*}(\{\underline{s}_i \in A_i\} \cap \{\underline{s}_{i+k} \in A_{i+k}\}) = \check{P}_{\pi_0^*}(\{\underline{s}_i \in A_i\})\check{P}(\{\underline{s}_{i+k} \in A_{i+k}\}).$$

Let  $0 \leq i_1 < i_2 < \dots < i_l$ . By the same pre- and post- process reasoning we obtain for  $A_{i_1}, \dots, A_{i_l}$  Borel subsets of  $R$  that

$$\begin{aligned} \check{P}_{\pi_0^*}(\{\underline{s}_{i_1} \in A_{i_1}\} \cap \dots \cap \{\underline{s}_{i_l} \in A_{i_l}\}) &= \\ &= \check{P}_{\pi_0^*}(\{\underline{s}_{i_1} \in A_{i_1}\} \cap \dots \cap \{\underline{s}_{i_{l-1}} \in A_{i_{l-1}}\}) \cdot \check{P}_{\pi_0^*}(\{\underline{s}_{i_l} \in A_{i_l}\}), \end{aligned}$$

and the proof is complete by induction.  $\square$

Now we turn to prove the following lemma, which generalizes the conclusions drawn in [Hobert & Robert 2004] for uniformly ergodic Markov chains.

**Lemma 3.2.2.** *Let the Markov chain  $(X_n)_{n \geq 0}$  be recurrent (and  $(X_{nm})_{n \geq 0}$  be recurrent) and let the minorization condition (2.9) hold with  $\pi(C) > 0$ . Then*

$$\mathcal{L}(X_{\tau_{\bar{\alpha}}(1)} | \{X_0, Y_0\} \in \bar{\alpha}) = \mathcal{L}(X_{\sigma_{\bar{\alpha}}(0)} | \{X_0, Y_0\} \sim \nu_m^*) = \pi_C(\cdot), \quad (3.6)$$

where  $\pi_C(\cdot)$  is a probability measure proportional to  $\pi$  truncated to  $C$ , that is  $\pi_C(B) = \pi(C)^{-1}\pi(B \cap C)$ .

*Proof.* The first equation in (3.6) is a straightforward consequence of the split chain construction. To prove the second one we use Theorem 2.2.8 for the

split  $m$ -skeleton with  $A = \check{\alpha}$ . Thus  $\tau_A = \tau_{\check{\alpha}}(1)$  and  $\check{\pi} := \pi^*$  is the invariant measure for the split  $m$ -skeleton. Let  $C \supseteq B \in \mathcal{B}(\mathcal{X})$ , and compute

$$\begin{aligned} \varepsilon\pi(B) &= \check{\pi}(B \times \{1\}) = \int_{\check{\alpha}} \check{E}_{x,y} \left[ \sum_{n=1}^{\tau_{\check{\alpha}}(1)} \mathbb{I}_{B \times \{1\}}(X_{nm}, Y_n) \right] \check{\pi}(dx, dy) \\ &= \check{\pi}(\check{\alpha}) \check{E}_{\nu_m^*} \left[ \sum_{n=0}^{\sigma_{\check{\alpha}}(0)} \mathbb{I}_{B \times \{1\}}(X_{nm}, Y_n) \right] = \check{\pi}(\check{\alpha}) \check{E}_{\nu_m^*} \mathbb{I}_B(X_{\sigma_{\check{\alpha}}(0)}). \end{aligned}$$

This implies proportionality and the proof is complete.  $\square$

**Lemma 3.2.3.**  $\check{E}_{\pi_0^*} \bar{s}_i^2 \leq \frac{m^2 \pi \bar{g}^2}{\varepsilon \pi(C)} < \infty$  and  $(\bar{s}_i)_{i \geq 0}$  are 1-dependent identically distributed r.v.'s.

*Proof.* Recall that  $\bar{s}_i = \sum_{k=0}^{m-1} \bar{g}(X_{\sigma_{\check{\alpha}}(i+1)m+k}) - \check{E}_{\pi_0^*} \left( \sum_{k=0}^{m-1} \bar{g}(X_{\sigma_{\check{\alpha}}(i+1)m+k}) \right)$  and is a function of the random variable

$$\{X_{\sigma_{\check{\alpha}}(i+1)m}, \dots, X_{\sigma_{\check{\alpha}}(i+1)m+m-1}\}. \quad (3.7)$$

By  $\mu_i(\cdot)$  denote the distribution of (3.7) on  $\mathcal{X}^m$ . We will show that  $\mu_i$  does not depend on  $i$ . From (2.13), (2.17) and the Bayes rule, for  $x \in C$ , we obtain

$$\check{P}\left(X_{nm+1} \in dx_1, \dots, X_{(n+1)m-1} \in dx_{m-1}, X_{(n+1)m} \in dy\right) \quad (3.8)$$

$$\left| Y_0^n, X_0^{nm}; Y_n = 1, X_{nm} = x \right) = \frac{\nu_m(dy)}{P^m(x, dy)} P(x, dx_1) \cdots P(x_{m-1}, dy).$$

Lemma 3.2.2 together with (3.8) yields

$$\check{P}\left(X_{nm} \in dx, X_{nm+1} \in dx_1, \dots, X_{(n+1)m-1} \in dx_{m-1}, X_{(n+1)m} \in dy\right) \quad (3.9)$$

$$\left| Y_0^n, X_0^{nm}; Y_n = 1; \sigma_{\check{\alpha}}(0) < n \right) = \pi_C(dx) \frac{\nu_m(dy)}{P^m(x, dy)} P(x, dx_1) \cdots P(x_{m-1}, dy).$$

Note that  $\frac{\nu_m(dy)}{P^m(x, dy)}$  is just a Radon-Nykodym derivative and thus (3.9) is a well defined measure on  $\mathcal{X}^{m+1}$ , say  $\mu(\cdot)$ . It remains to notice, that  $\mu_i(A) = \mu(A \times \mathcal{X})$  for any Borel  $A \subset \mathcal{X}^m$ . Thus  $\mu_i$ ,  $i \geq 0$  are identical and hence  $\bar{s}_i$ ,  $i \geq 0$  have the same distribution. Due to Lemma 2.2.9 we obtain that  $\bar{s}_i$ ,



$i \geq 0$  are 1-dependent. To prove  $\check{E}_{\pi_0^*} \check{s}_i^2 < \infty$ , we first note that  $\frac{\nu_m(dy)}{P^m(x, dy)} \leq 1/\varepsilon$  and also  $\pi_C(\cdot) \leq \frac{1}{\pi(C)}\pi(\cdot)$ . Hence

$$\mu_i(A) = \mu(A \times \mathcal{X}) \leq \frac{1}{\varepsilon\pi(C)}\mu_{\text{chain}}(A),$$

where  $\mu_{\text{chain}}$  is defined by  $\pi(dx)P(x, dx_1) \dots P(x_{m-2}, dx_{m-1})$ . Thus

$$\left| \check{E}_{\pi_0^*} \left( \sum_{k=0}^{m-1} \bar{g}(X_{\sigma_{\check{\alpha}}(i+1)m+k}) \right) \right| \leq \frac{m\pi|\bar{g}|}{\varepsilon\pi(C)} < \infty.$$

Now let  $\check{s}_i = \sum_{k=0}^{m-1} \bar{g}(X_{\sigma_{\check{\alpha}}(i+1)m+k})$  and proceed

$$\begin{aligned} \check{E}_{\pi_0^*} \check{s}_i^2 &\leq \check{E}_{\pi_0^*} \check{s}_i^2 \leq \frac{1}{\varepsilon\pi(C)}\mu_{\text{chain}} \check{s}_i^2 = \frac{1}{\varepsilon\pi(C)} E_{\pi} \left( \sum_{k=0}^{m-1} \bar{g}(X_k) \right)^2 \\ &\leq \frac{m}{\varepsilon\pi(C)} E_{\pi} \left[ \sum_{k=0}^{m-1} \bar{g}^2(X_k) \right] \leq \frac{m^2\pi\bar{g}^2}{\varepsilon\pi(C)}. \end{aligned}$$

□

We need a result which gives the connection between stochastic boundedness and the existence of the second moment of  $\underline{s}_i$ . We state it in a general form.

**Theorem 3.2.4.** *Let  $(X_n)_{n \geq 0}$  be a sequence of independent identically distributed random variables and  $S_n = \sum_{k=0}^{n-1} X_k$ . Suppose that  $(\tau_n)$  is a sequence of positive, integer valued r.v.'s such that  $\tau_n/n \rightarrow a \in (0, \infty)$  in probability when  $n \rightarrow \infty$  and the sequence  $(n^{-1/2}S_{\tau_n})$  is stochastically bounded. Then  $EX_0^2 < \infty$  and  $EX_0 = 0$ .*

The proof of Theorem 3.2.4 is based on the following lemmas.

**Lemma 3.2.5.** *Let  $\delta \in (0, 1)$  and  $t_0 := \sup\{t > 0: \sup_{0 \leq k \leq n} P(|S_k| \geq t) \geq \delta\}$ . Then  $P(|S_{10n}| \geq 4t_0) \geq (1-\delta)(\delta/4)^{20}$  and  $P(\sup_{k \leq n} |S_k| \leq 3t_0) \geq 1-3\delta$ .*

*Proof.* By the definition of  $t_0$  there exists  $0 \leq n_0 \leq n$  such that  $P(|S_{n_0}| \geq t_0) \geq \delta$ . Then either  $P(|S_n| \geq t_0/2) \geq \delta/2$  or  $P(|S_n| \geq t_0/2) < \delta/2$  and consequently

$$\begin{aligned} P(|S_{n-n_0}| \geq t_0/2) &= P(|S_n - S_{n_0}| \geq t_0/2) \\ &\geq P(|S_{n_0}| \geq t_0) - P(|S_n| \geq t_0/2) \geq \delta/2. \end{aligned}$$

Thus there exists  $n/2 \leq n_1 \leq n$  such that  $P(|S_{n_1}| \geq t_0/2) \geq \delta/2$ . Let  $10n = an_1 + b$  with  $0 \leq b < n_1$ , then  $10 \leq a \leq 20$ ,

$$\begin{aligned} P(|S_{an_1}| \geq 5t_0) &\geq P(S_{an_1} \geq at_0/2) + P(S_{an_1} \leq -at_0/2) \\ &\geq (P(S_{n_1} \geq t_0/2))^a + (P(S_{n_1} \leq -t_0/2))^a \geq (\delta/4)^a, \end{aligned}$$

hence

$$\begin{aligned} P(|S_{10n}| \geq 4t_0) &\geq P(|S_{an_1}| \geq 5t_0)P(|S_{10n} - S_{an_1}| \leq t_0) \\ &\geq (\delta/4)^a(1 - \delta) \geq (1 - \delta)(\delta/4)^{20}. \end{aligned}$$

Finally by the Levy-Octaviani inequality we obtain

$$P\left(\sup_{k \leq n} |S_k| > 3t_0\right) \leq 3 \sup_{k \leq n} P(|S_k| > t_0) \leq 3\delta.$$

□

**Lemma 3.2.6.** *Let  $c^2 < \text{Var}(X_1)$ , then for sufficiently large  $n$ ,  $P(|S_n| \geq c\sqrt{n}/4) \geq 1/16$ .*

*Proof.* Let  $(X'_i)$  be an independent copy of  $(X_i)$  and  $S'_k = \sum_{i=1}^k X'_i$ . Moreover let  $(\varepsilon_i)$  be a sequence of independent symmetric  $\pm 1$  r.v.'s, independent of  $(X_i)$  and  $(X'_i)$ . For any reals  $(a_i)$  we get by the Paley-Zygmund inequality,

$$\begin{aligned} P\left(\left|\sum_{i=1}^n a_i \varepsilon_i\right| \geq \frac{1}{2} \left(\sum_i a_i^2\right)^{1/2}\right) &= P\left(\left|\sum_{i=1}^n a_i \varepsilon_i\right|^2 \geq \frac{1}{4} E\left|\sum_{i=1}^n a_i \varepsilon_i\right|^2\right) \\ &\geq \left(1 - \frac{1}{4}\right)^2 \frac{(E|\sum_{i=1}^n a_i \varepsilon_i|^2)^2}{E|\sum_{i=1}^n a_i \varepsilon_i|^4} \geq \frac{3}{16}. \end{aligned}$$

Hence

$$\begin{aligned} P\left(|S_n - S'_n| \geq \frac{c}{2}\sqrt{n}\right) &= P\left(\left|\sum_{i=1}^n \varepsilon_i(X_i - X'_i)\right| \geq \frac{c}{2}\sqrt{n}\right) \\ &\geq \frac{3}{16} P\left(\sum_{i=1}^n (X_i - X'_i)^2 \geq c^2 n\right) \geq \frac{1}{8} \end{aligned}$$

for sufficiently large  $n$  by the Weak LLN. Thus

$$\begin{aligned} \frac{1}{8} \leq P\left(|S_n - S'_n| \geq \frac{c}{2}\sqrt{n}\right) &\leq P\left(|S_n| \geq \frac{c}{4}\sqrt{n}\right) + P\left(|S'_n| \geq \frac{c}{4}\sqrt{n}\right) \\ &\leq 2P\left(|S_n| \geq \frac{c}{4}\sqrt{n}\right). \end{aligned}$$

□

**Corollary 3.2.7.** *Let  $c^2 < \text{Var}(X_1)$ , then for sufficiently large  $n$ ,*

$$P\left(\inf_{10n \leq k \leq 11n} |S_k| \geq \frac{1}{4}c\sqrt{n}\right) \geq 2^{-121}.$$

*Proof.* Let  $t_0$  be as in Lemma 3.2.5 for  $\delta = 1/16$ , then

$$\begin{aligned} P\left(\inf_{10n \leq k \leq 11n} |S_k| \geq t_0\right) &\geq P\left(|S_{10n}| \geq 4t_0, \sup_{10n \leq k \leq 11n} |S_k - S_{10n}| \leq 3t_0\right) \\ &= P(|S_{10n}| \geq 4t_0)P\left(\sup_{k \leq n} |S_k| \leq 3t_0\right) \geq 2^{-121}. \end{aligned}$$

Hence by Lemma 3.2.5 we obtain  $t_0 \geq c\sqrt{n}/4$  for large  $n$ .  $\square$

*Proof of Theorem 3.2.4.* By Corollary 3.2.7 for any  $c^2 < \text{Var}(X)$  we have,

$$\begin{aligned} P\left(|S_{\tau_n}| \geq \frac{c}{20}\sqrt{an}\right) &\geq P\left(\left|\frac{\tau_n}{n} - a\right| \leq \frac{a}{21}, \inf_{\frac{20}{21}an \leq k \leq \frac{22}{21}an} |S_k| \geq \frac{c}{20}\sqrt{an}\right) \geq \\ &\geq P\left(\inf_{\frac{20}{21}an \leq k \leq \frac{22}{21}an} |S_k| \geq \frac{c}{4}\sqrt{\frac{2an}{21}}\right) - P\left(\left|\frac{\tau_n}{n} - a\right| > \frac{a}{21}\right) \\ &\geq 2^{-121} - P\left(\left|\frac{\tau_n}{n} - a\right| > \frac{a}{21}\right) \geq 2^{-122} \end{aligned}$$

for sufficiently large  $n$ . Since  $(n^{-1/2}S_{\tau_n})$  is stochastically bounded, we immediately obtain  $\text{Var}(X_1) < \infty$ . If  $EX_1 \neq 0$  then

$$\left|\frac{1}{\sqrt{n}}S_{\tau_n}\right| = \left|\frac{S_{\tau_n}}{\tau_n}\right| \left|\frac{\tau_n}{n}\right| \sqrt{n} \rightarrow \infty \quad \text{in probability when } n \rightarrow \infty.$$

$\square$

### 3.3 A Characterization of $\sqrt{n}$ -CLTs

In this section we provide a generalization of Theorem 17.3.6 of [Meyn & Tweedie 1993]. We obtain an if and only if condition for the  $\sqrt{n}$ -CLT in terms of finiteness of the second moment of a centered excursion from  $\check{\alpha}$ .

**Theorem 3.3.1.** *Suppose that  $(X_n)_{n \geq 0}$  is ergodic and  $\pi(g^2) < \infty$ . Let  $\nu_m$  be the measure satisfying (2.9), then the  $\sqrt{n}$ -CLT holds if and only if*

$$\check{E}_{\nu_m^*} \left[ \left( \sum_{n=0}^{\sigma_{\check{\alpha}}(0)} Z_n(\check{g}) \right)^2 \right] < \infty. \quad (3.10)$$

Furthermore we have the following formula for variance

$$\sigma_g^2 = \frac{\varepsilon\pi(C)}{m} \left\{ \check{E}_{\nu_m^*} \left[ \left( \sum_{n=0}^{\sigma_{\check{\alpha}}(0)} Z_n(\bar{g}) \right)^2 \right] + 2\check{E}_{\nu_m^*} \left[ \left( \sum_{n=0}^{\sigma_{\check{\alpha}}(0)} Z_n(\bar{g}) \right) \left( \sum_{n=\sigma_{\check{\alpha}}(0)+1}^{\sigma_{\check{\alpha}}(1)} Z_n(\bar{g}) \right) \right] \right\}.$$

*Proof.* For  $n \geq 0$  define

$$l_n := \max\{k \geq 1 : m(\sigma_{\check{\alpha}}(k) + 1) \leq n\}$$

and for completeness  $l_n := 0$  if  $m(\sigma_{\check{\alpha}}(0) + 1) \geq n$ . First we are going to show that

$$\left| \frac{1}{\sqrt{n}} \sum_{j=0}^{n-1} \bar{g}(X_j) - \frac{1}{\sqrt{n}} \sum_{j=0}^{l_n-1} s_j \right| \rightarrow 0 \quad \text{in probability.} \quad (3.11)$$

Thus we have to verify that the initial and final terms of the sum do not matter. First observe that by the Harris recurrence property of the chain  $\sigma_{\check{\alpha}}(0) < \infty$ ,  $\check{P}_{\pi_0^*}$ -a.s. and hence  $\lim_{n \rightarrow \infty} \check{P}_{\pi_0^*}(m\sigma_{\check{\alpha}}(0) \geq n) = 0$  and  $\check{P}_{\pi_0^*}(\sigma_{\check{\alpha}}(0) < \infty) = 1$ . This yields

$$\left| \frac{1}{\sqrt{n}} \sum_{j=0}^{n-1} \bar{g}(X_j) - \frac{1}{\sqrt{n}} \sum_{j=m(\sigma_{\check{\alpha}}(0)+1)}^{n-1} \bar{g}(X_j) \right| \rightarrow 0, \quad \check{P} - \text{a.s.} \quad (3.12)$$

The second point is to provide a similar argument for the tail terms and to show that

$$\left| \frac{1}{\sqrt{n}} \sum_{j=m(\sigma_{\check{\alpha}}(0)+1)}^{n-1} \bar{g}(X_j) - \frac{1}{\sqrt{n}} \sum_{j=m(\sigma_{\check{\alpha}}(0)+1)}^{m\sigma_{\check{\alpha}}(l_n)+m-1} \bar{g}(X_j) \right| \rightarrow 0, \quad \text{in probability.} \quad (3.13)$$

For  $\varepsilon > 0$  we have

$$\begin{aligned} \check{P}_{\pi_0^*} \left( \left| \frac{1}{\sqrt{n}} \sum_{j=m(\sigma_{\check{\alpha}}(l_n)+1)}^{n-1} \bar{g}(X_j) \right| > \varepsilon \right) &\leq \check{P}_{\pi_0^*} \left( \frac{1}{\sqrt{n}} \sum_{j=\sigma_{\check{\alpha}}(l_n)+1}^{\sigma_{\check{\alpha}}(l_n)+1} Z_j(|\bar{g}|) > \varepsilon \right) \\ &\leq \sum_{k=0}^{\infty} \check{P}_{\check{\alpha}} \left( \frac{1}{\sqrt{n}} \sum_{j=1}^{\tau_{\check{\alpha}}(1)} Z_j(|\bar{g}|) > \varepsilon, \tau_{\check{\alpha}}(1) \geq k \right). \end{aligned}$$

Now since  $\sum_{k=0}^{\infty} \check{P}_{\check{\alpha}}(\tau_{\check{\alpha}}(1) \geq k) \leq \check{E}_{\check{\alpha}}\tau_{\check{\alpha}}(1) < \infty$ , where we use that  $\check{\alpha}$  is an atom for the split chain, we deduce from the Lebesgue majorized convergence theorem that (3.13) holds. Obviously (3.12) and (3.13) yield (3.11).

We turn to prove that the condition (3.10) is sufficient for the CLT to hold. We will show that random numbers  $l_n$  can be replaced by their non-random equivalents. Namely we apply the LLN (Theorem 17.3.2 in [Meyn & Tweedie 1993]) to ensure that

$$\lim_{n \rightarrow \infty} \frac{l_n}{n} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^{\lfloor n/m \rfloor - 1} \mathbb{I}_{\{(X_{mk}, Y_k) \in \check{\alpha}\}} = \frac{\check{\pi}(\check{\alpha})}{m}, \quad \check{P}_{\pi_0^*} - \text{a.s.} \quad (3.14)$$

Let

$$n^* := \lfloor \check{\pi}(\check{\alpha})nm^{-1} \rfloor, \quad \underline{n} := \lceil (1-\varepsilon)\check{\pi}(\check{\alpha})nm^{-1} \rceil, \quad \bar{n} := \lfloor (1+\varepsilon)\check{\pi}(\check{\alpha})nm^{-1} \rfloor.$$

Due to the LLN we know that for any  $\varepsilon > 0$ , there exists  $n_0$  such that for all  $n \geq n_0$  we have  $\check{P}_{\pi_0^*}(\underline{n} \leq l_n \leq \bar{n}) \geq 1 - \varepsilon$ . Consequently

$$\begin{aligned} \check{P}_{\pi_0^*} \left( \left| \sum_{j=0}^{l_n-1} s_j - \sum_{j=0}^{n^*} s_j \right| > \sqrt{n}\beta \right) &\leq \varepsilon + \check{P}_{\pi_0^*} \left( \max_{\underline{n} \leq l \leq n^*} \left| \sum_{j=l}^{n^*} s_j \right| > \beta\sqrt{n} \right) + \\ &+ \check{P}_{\pi_0^*} \left( \max_{n^*+1 \leq l \leq \bar{n}} \left| \sum_{j=n^*+1}^l s_j \right| > \beta\sqrt{n} \right). \end{aligned} \quad (3.15)$$

Since  $(s_j)_{j \geq 0}$  are 1-dependent,  $M_k := \sum_{j=0}^k s_j$  is not necessarily a martingale. Thus to apply the classical Kolmogorov inequality we define  $M_k^0 = \sum_{j=0}^{\infty} s_{2j} \mathbb{I}_{\{2j \leq k\}}$  and  $M_k^1 = \sum_{j=0}^{\infty} s_{1+2j} \mathbb{I}_{\{1+2j \leq k\}}$ , which are clearly square-integrable martingales (due to (3.10)). Hence

$$\begin{aligned} \check{P}_{\pi_0^*} \left( \max_{\underline{n} \leq l \leq n^*} |M_{n^*} - M_l| > \beta\sqrt{n} \right) &\leq \check{P}_{\pi_0^*} \left( \max_{\underline{n} \leq l \leq n^*} |M_{n^*}^0 - M_l^0| > \frac{\beta\sqrt{n}}{2} \right) + \\ &+ \check{P}_{\pi_0^*} \left( \max_{\underline{n} \leq l \leq n^*} |M_{n^*}^1 - M_l^1| > \frac{\beta\sqrt{n}}{2} \right) \\ &\leq \frac{4}{n\beta^2} \sum_{k=0}^1 (\check{E}_{\pi_0^*} |M_{n^*}^k - M_{\underline{n}}^k|^2) \\ &\leq C\varepsilon\beta^{-2} \check{E}_{\nu_m^*}(s_0^2), \end{aligned} \quad (3.16)$$

where  $C$  is a universal constant. In the same way we show that

$$\check{P} \left( \max_{n^*+1 \leq l \leq \bar{n}} |M_l - M_{n^*+1}| > \beta\sqrt{n} \right) \leq C\varepsilon\beta^{-2} \check{E}_{\nu_m^*}(s_0^2),$$

consequently, since  $\varepsilon$  is arbitrary, we obtain

$$\left| \frac{1}{\sqrt{n}} \sum_{j=0}^{l_n-1} s_j - \frac{1}{\sqrt{n}} \sum_{j=0}^{n^*} s_j \right| \rightarrow 0, \quad \text{in probability.} \quad (3.17)$$

The last step is to provide an argument for the CLT for 1-dependent, identically distributed random variables. Namely, we have to prove that

$$\frac{1}{\sqrt{n}} \sum_{j=0}^n s_j \xrightarrow{d} \mathcal{N}(0, \bar{\sigma}^2), \quad \text{as } n \rightarrow \infty, \quad (3.18)$$

where

$$\bar{\sigma}^2 := \check{E}_{\nu_m^*}(s_0(\bar{g}))^2 + 2\check{E}_{\nu_m^*}(s_0(\bar{g})s_1(\bar{g})).$$

Observe that (3.12), (3.13), (3.17) and (3.18) imply Theorem 3.3.1. We fix  $k \geq 2$  and define  $\xi_j := s_{kj+1}(\bar{g}) + \dots + s_{kj+k-1}(\bar{g})$ , consequently  $\xi_j$  are i.i.d. random variables and

$$\frac{1}{\sqrt{n}} \sum_{j=0}^n s_j = \frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor n/k \rfloor - 1} \xi_j + \frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor n/k \rfloor} s_{kj}(\bar{g}) + \frac{1}{\sqrt{n}} \sum_{j=k\lfloor n/k \rfloor + 1}^n s_j. \quad (3.19)$$

Obviously the last term converges to 0 in probability. Denoting

$$\begin{aligned} \sigma_k^2 &:= \check{E}_{\pi_0^*}(\xi_j)^2 = (k-1)\check{E}_{\nu_m^*}(s_0(\bar{g}))^2 + 2(k-2)\check{E}_{\nu_m^*}(s_0(\bar{g})s_1(\bar{g})), \\ \sigma_s^2 &:= \check{E}_{\nu_m^*}(s_0(\bar{g}))^2. \end{aligned}$$

we use the classical CLT for i.i.d. random variables to see that

$$\frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor n/k \rfloor - 1} \xi_j \xrightarrow{d} \mathcal{N}(0, k^{-1}\sigma_k^2), \quad \text{and} \quad \frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor n/k \rfloor} s_{kj}(\bar{g}) \xrightarrow{d} \mathcal{N}(0, k^{-1}\sigma_s^2). \quad (3.20)$$

Moreover

$$\lim_{n \rightarrow \infty} \left[ \frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor n/k \rfloor - 1} \xi_j + \frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor n/k \rfloor} s_{kj}(\bar{g}) \right] \quad (3.21)$$

converges to  $\mathcal{N}(0, \sigma_g^2)$ , with  $k \rightarrow \infty$ . Since the weak convergence is metrizable we deduce from (3.19), (3.20) and (3.21) that (3.18) holds.

The remaining part is to prove that (3.10) is also necessary for the CLT to hold. Note that if  $\sum_{k=0}^n \bar{g}(X_k)/\sqrt{n}$  verifies the CLT then  $\sum_{j=0}^{l_n-1} s_j$  is stochastically bounded by (3.11). We use the decomposition  $s_i = \underline{s}_i + \bar{s}_i$ ,  $i \geq 0$  introduced in Section 3.2. By Lemma 3.2.3 we know that  $\bar{s}_j$  is a sequence of 1-dependent random variables with the same distribution and finite second moment. Thus from the first part of the proof we deduce that  $\sum_{j=0}^{l_n-1} \bar{s}_j/\sqrt{n}$  verifies a CLT and thus is stochastically bounded. Consequently the remaining sequence  $\sum_{j=0}^{l_n-1} \underline{s}_j/\sqrt{n}$  also must be stochastically bounded. Lemma 3.2.1 states that  $(\underline{s}_j)_{j \geq 0}$  is a sequence of i.i.d. r.v.'s, hence  $\check{E}[\underline{s}_j^2] < \infty$  by Theorem 3.2.4. Also  $l_n/n \rightarrow \check{\pi}(\check{\alpha})m^{-1}$  by (3.14). Applying the inequality  $(a+b)^2 \leq 2(a^2 + b^2)$  we obtain

$$\check{E}_{\pi_0^*}[s_j]^2 \leq 2(\check{E}_{\pi_0^*}[\underline{s}_j^2] + \check{E}_{\pi_0^*}[\bar{s}_j^2]) < \infty$$

which completes the proof.  $\square$

*Remark 3.3.2.* Note that in the case of  $m = 1$  we have  $\bar{s}_i \equiv 0$  and for Theorem 3.3.1 to hold, it is enough to assume  $\pi|g| < \infty$  instead of  $\pi(g^2) < \infty$ . In the case of  $m > 1$ , assuming only  $\pi|g| < \infty$  and (3.10) implies the  $\sqrt{n}$ -CLT, but the proof of the converse statement fails, and in fact the converse statement does not hold (one can easily provide an appropriate counterexample).

## 3.4 Uniform Ergodicity

In view of Theorem 3.3.1 providing a regeneration proof of Theorem 3.1.4 amounts to establishing conditions (3.10) and checking the formula for the asymptotic variance. To this end we need some additional facts about small sets for uniformly ergodic Markov chains.

**Theorem 3.4.1.** *If  $(X_n)_{n \geq 0}$ , a Markov chain on  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  with stationary distribution  $\pi$  is uniformly ergodic, then  $\mathcal{X}$  is  $\nu_m$ -small for some  $\nu_m$ .*

Hence for uniformly ergodic chains (2.9) holds for all  $x \in \mathcal{X}$ . Theorem 3.4.1 is well known in literature, in particular it results from Theorems 5.2.1 and 5.2.4 in [Meyn & Tweedie 1993] with their  $\psi = \pi$ .

Theorem 3.4.1 implies that for uniformly ergodic Markov chains (2.10) can be rewritten as

$$P^m(x, \cdot) = \varepsilon \nu_m(\cdot) + (1 - \varepsilon)R(x, \cdot). \quad (3.22)$$

The following mixture representation of  $\pi$  will turn out very useful.

**Lemma 3.4.2.** *If  $(X_n)_{n \geq 0}$  is an ergodic Markov chain with transition kernel  $P$  and (3.22) holds, then*

$$\pi = \varepsilon \mu := \varepsilon \sum_{n=0}^{\infty} \nu_m (1 - \varepsilon)^n R^n. \quad (3.23)$$

*Remark 3.4.3.* This can be easily extended to the more general setting than this of uniformly ergodic chains, namely let  $P^m(x, \cdot) = s(x)\nu_m(\cdot) + (1 - s(x))R(x, \cdot)$ ,  $s : \mathcal{X} \rightarrow [0, 1]$ ,  $\pi s > 0$ . In this case  $\pi = \pi s \sum_{n=0}^{\infty} \nu_m R_{\#}^n$ , where  $R_{\#}(x, \cdot) = (1 - s(x))R(x, \cdot)$ . Related decompositions under various assumptions can be found e.g. in [Nummelin 2002], [Hobert & Robert 2004] and [Breyer & Roberts 2001] and are closely related to perfect sampling algorithms, such as coupling from the past (CFTP) introduced in [Propp & Wilson 1996].

*Proof.* First check that the measure in question is a probability measure.

$$\left( \varepsilon \sum_{n=0}^{\infty} \nu_m (1 - \varepsilon)^n R^n \right) (\mathcal{X}) = \varepsilon \sum_{n=0}^{\infty} (1 - \varepsilon)^n (\nu_m R^n) (\mathcal{X}) = 1.$$

It is also invariant for  $P^m$  :

$$\begin{aligned} \left( \sum_{n=0}^{\infty} \nu_m (1 - \varepsilon)^n R^n \right) P^m &= \left( \sum_{n=0}^{\infty} \nu_m (1 - \varepsilon)^n R^n \right) (\varepsilon \nu_m + (1 - \varepsilon)R) \\ &= \varepsilon \mu \nu_m + \sum_{n=1}^{\infty} \nu_m (1 - \varepsilon)^n R^n = \sum_{n=0}^{\infty} \nu_m (1 - \varepsilon)^n R^n. \end{aligned}$$

Hence by ergodicity  $\varepsilon \mu = \varepsilon \mu P^{nm} \rightarrow \pi$ , as  $n \rightarrow \infty$ . This completes the proof.  $\square$

**Corollary 3.4.4.** *The decomposition in Lemma 3.4.2 implies that*

$$\begin{aligned} (i) \quad \check{E}_{\nu_m^*} \left( \sum_{n=0}^{\sigma(0)} \mathbb{I}_{\{X_{nm} \in A\}} \right) &= \check{E}_{\nu_m^*} \left( \sum_{n=0}^{\infty} \mathbb{I}_{\{X_{nm} \in A\}} \mathbb{I}_{\{Y_0=0, \dots, Y_{n-1}=0\}} \right) = \varepsilon^{-1} \pi(A), \\ (ii) \quad \check{E}_{\nu_m^*} \left( \sum_{n=0}^{\infty} f(X_{nm}, X_{nm+1}, \dots; Y_n, Y_{n+1}, \dots) \mathbb{I}_{\{Y_0=0, \dots, Y_{n-1}=0\}} \right) &= \\ &= \varepsilon^{-1} \check{E}_{\pi^*} f(X_0, X_1, \dots; Y_0, Y_1, \dots). \end{aligned}$$



*Proof.* (i) is a direct consequence of (3.23). To see (ii) note that  $Y_n$  is a coin toss independent of  $\{Y_0, \dots, Y_{n-1}\}$  and  $X_{nm}$ , this allows for  $\pi^*$  instead of  $\pi$  on the RHS of (ii). Moreover the evolution of  $\{X_{nm+1}, X_{nm+2}, \dots; Y_{n+1}, Y_{n+2}, \dots\}$  depends only (and explicitly by (2.13) and (2.14)) on  $X_{nm}$  and  $Y_n$ . Now use (i).  $\square$

Our object of interest is

$$\begin{aligned}
I &= \check{E}_{\nu_m^*} \left[ \left( \sum_{n=0}^{\sigma(0)} Z_n(\bar{g}) \right)^2 \right] = \check{E}_{\nu_m^*} \left[ \left( \sum_{n=0}^{\infty} Z_n(\bar{g}) \mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq n\}} \right)^2 \right] \\
&= \check{E}_{\nu_m^*} \left[ \sum_{n=0}^{\infty} Z_n(\bar{g})^2 \mathbb{I}_{\{Y_0=0, \dots, Y_{n-1}=0\}} \right] + \\
&\quad + 2\check{E}_{\nu_m^*} \left[ \sum_{n=0}^{\infty} \sum_{k=n+1}^{\infty} Z_n(\bar{g}) \mathbb{I}_{\{\sigma(0) \geq n\}} Z_k(\bar{g}) \mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq k\}} \right] \\
&= A + B \tag{3.24}
\end{aligned}$$

Next we use Corollary 3.4.4 and then the inequality  $2ab \leq a^2 + b^2$  to bound the term  $A$  in (3.24).

$$\begin{aligned}
A &= \varepsilon^{-1} \check{E}_{\pi^*} Z_0(\bar{g})^2 = \varepsilon^{-1} E_{\pi} \left( \sum_{k=0}^{m-1} \bar{g}(X_k) \right)^2 \\
&\leq \varepsilon^{-1} m E_{\pi} \left[ \sum_{k=0}^{m-1} \bar{g}^2(X_k) \right] \leq \varepsilon^{-1} m^2 \pi \bar{g}^2 < \infty.
\end{aligned}$$

We proceed similarly with the term  $B$

$$\begin{aligned}
|B| &\leq 2\check{E}_{\nu_m^*} \left[ \sum_{n=0}^{\infty} |Z_n(\bar{g})| \mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq n\}} \sum_{k=1}^{\infty} |Z_{n+k}(\bar{g})| \mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq n+k\}} \right] \\
&= 2\varepsilon^{-1} \check{E}_{\pi^*} \left[ |Z_0(\bar{g})| \sum_{k=1}^{\infty} |Z_k(\bar{g})| \mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq k\}} \right].
\end{aligned}$$

By Cauchy-Schwarz,

$$\begin{aligned}
\check{E}_{\pi^*} [\mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq k\}} |Z_0(\bar{g})| |Z_k(\bar{g})|] &\leq \sqrt{\check{E}_{\pi^*} [\mathbb{I}_{\{\sigma_{\bar{\alpha}}(0) \geq k\}} Z_0(\bar{g})^2]} \sqrt{\check{E}_{\pi^*} Z_k(\bar{g})^2} \\
&= \sqrt{\check{E}_{\pi^*} [\mathbb{I}_{\{Y_0=0\}} \mathbb{I}_{\{Y_1=0, \dots, Y_{k-1}=0\}} Z_0(\bar{g})^2]} \sqrt{\check{E}_{\pi^*} Z_0(\bar{g})^2}.
\end{aligned}$$

Observe that  $\{Y_1, \dots, Y_{k-1}\}$  and  $\{X_0, \dots, X_{m-1}\}$  are independent. We drop  $\mathbb{I}_{\{Y_0=0\}}$  to obtain

$$\check{E}_{\pi^*} [\mathbb{I}_{\{\sigma_{\check{\alpha}}(0) \geq k\}} |Z_0(\bar{g})| |Z_k(\bar{g})|] \leq (1 - \varepsilon)^{\frac{k-1}{2}} \check{E}_{\pi^*} Z_0(\bar{g})^2 \leq (1 - \varepsilon)^{\frac{k-1}{2}} m^2 \pi g^2.$$

Hence  $|B| < \infty$ , and the proof of (3.10) is complete. To get the variance formula note that the convergence we have established implies

$$I = \varepsilon^{-1} \check{E}_{\pi^*} \left[ Z_0(\bar{g}) \right]^2 + 2\varepsilon^{-1} \check{E}_{\pi^*} \left[ Z_0(\bar{g}) \sum_{k=1}^{\infty} Z_k(\bar{g}) \mathbb{I}_{\{\sigma_{\check{\alpha}}(0) \geq k\}} \right].$$

Similarly we obtain

$$\begin{aligned} J &:= 2\check{E}_{\nu_m^*} \left[ \left( \sum_{n=0}^{\sigma_{\check{\alpha}}(0)} Z_n(\bar{g}) \right) \left( \sum_{n=\sigma_{\check{\alpha}}(0)+1}^{\sigma_{\check{\alpha}}(1)} Z_n(\bar{g}) \right) \right] \\ &= 2\varepsilon^{-1} \check{E}_{\pi^*} \left[ Z_0(\bar{g}) \sum_{k=\sigma_{\check{\alpha}}(0)+1}^{\infty} Z_k(\bar{g}) \mathbb{I}_{\{\sigma_{\check{\alpha}}(1) \geq k\}} \right]. \end{aligned}$$

Since  $\pi(C) = 1$ , we have  $\sigma_g^2 = \varepsilon m^{-1} (I + J)$ . Next we use Lemma 2.2.9 and  $\check{E}_{\pi^*} Z_0(\bar{g}) = 0$  to drop indicators and since for  $f : \mathcal{X} \rightarrow R$ , also  $\check{E}_{\pi^*} f = E_{\pi} f$ , we have

$$\begin{aligned} \varepsilon(I + J) &= \check{E}_{\pi^*} \left[ Z_0(\bar{g}) \left( Z_0(\bar{g}) + 2 \sum_{k=1}^{\infty} Z_k(\bar{g}) \right) \right] \\ &= E_{\pi} \left[ Z_0(\bar{g}) \left( Z_0(\bar{g}) + 2 \sum_{k=1}^{\infty} Z_k(\bar{g}) \right) \right]. \end{aligned}$$

Now, since all the integrals are taken with respect to the stationary measure, we can for a moment assume that the chain runs in stationarity from  $-\infty$  rather than starts at time 0 with  $X_0 \sim \pi$ . Thus

$$\begin{aligned} \sigma_g^2 &= m^{-1} E_{\pi} \left[ Z_0(\bar{g}) \left( \sum_{k=-\infty}^{\infty} Z_k(\bar{g}) \right) \right] = m^{-1} E_{\pi} \left[ \sum_{l=0}^{m-1} \bar{g}(X_l) \left( \sum_{k=-\infty}^{\infty} \bar{g}(X_k) \right) \right] \\ &= E_{\pi} \left[ \bar{g}(X_0) \sum_{k=-\infty}^{\infty} \bar{g}(X_k) \right] = \int_{\mathcal{X}} \bar{g}^2 d\pi + 2 \int_{\mathcal{X}} \sum_{n=1}^{\infty} \bar{g}(X_0) \bar{g}(X_n) d\pi. \end{aligned}$$

### 3.5 The difference between $m = 1$ and $m \neq 1$

Assume the small set condition (2.9) holds and consider the split chain defined by (2.13) and (2.14). The following tours

$$\{ \{ X_{(\sigma(n)+1)m}, X_{(\sigma(n)+1)m+1}, \dots, X_{(\sigma(n+1)+1)m-1} \}, n = 0, 1, \dots \}$$

that start whenever  $X_k \sim \nu_m$  are of crucial importance to the regeneration theory and are eagerly analyzed by researchers. In virtually every paper on the subject there is a claim these objects are independent identically distributed random variables. This claim is usually considered obvious and no proof is provided. However this is not true if  $m > 1$ . In fact formulas (2.13) and (2.14) should be convincing enough, as  $X_{mn+1}, \dots, X_{(n+1)m}$  given  $Y_n = 1$  and  $X_{nm} = x$  are linked in a way described by  $P(x, dx_1) \cdots P(x_{m-1}, dy)$ . In particular consider a Markov chain on  $\mathcal{X} = \{a, b, c, d, e\}$  with transition probabilities

$$\begin{aligned} P(a, b) &= P(a, c) = P(b, b) = P(b, d) = P(c, c) = P(c, e) = 1/2, \\ P(d, a) &= P(e, a) = 1. \end{aligned}$$

Let  $\nu_4(d) = \nu_4(e) = 1/2$  and  $\varepsilon = 1/8$ . Clearly  $P^4(x, \cdot) \geq \varepsilon \nu_4(\cdot)$  for every  $x \in \mathcal{X}$ , hence we established (2.9) with  $C = \mathcal{X}$ . Note that for this simplistic example each tour can start with  $d$  or  $e$ . However if it starts with  $d$  or  $e$  the previous tour must have ended with  $b$  or  $c$  respectively. This makes them dependent. Similar examples with general state space  $\mathcal{X}$  and  $C \neq \mathcal{X}$  can be easily provided. Hence Theorem 3.3.1 is critical to providing regeneration proofs of CLTs and standard arguments that involve i.i.d. random variables are not valid.

# Chapter 4

## Fixed-Width Asymptotics

Determining the length of simulation for MCMC algorithms that guarantees good quality of estimation is a fundamental problem. One possible approach is to wait until width of an asymptotic confidence interval based on the approximation by a normal distribution becomes smaller than a user-specified value. This requires estimating  $\sigma_g^2$  the variance of the asymptotic normal distribution. In this chapter we relax assumptions required to obtain strongly consistent estimators of  $\sigma_g^2$  in the regenerative setting.

Results of this chapter (in particular the key Lemma 4.3.3 and resulting from it Lemma 4.3.6 and Proposition 4.3.7) are based on the paper [Bednorz & Łatuszyński 2007] and are joint work with Witold Bednorz.

The presentation of the fixed-width asymptotic approach is based on [Jones et al. 2006]. We provide only a quick sketch, since the approach is well known in literature (see also [Geyer 1992], [Mykland et al. 1995], [Hobert et al. 2002]) and [Jones et al. 2006] is an excellent recent reference.

### 4.1 Asymptotic Confidence Intervals

Suppose that we are in the standard MCMC setting and our goal is to estimate  $I = E_\pi g = \int_{\mathcal{X}} g(x)\pi(dx)$ . Let  $(X_n)_{n \geq 0}$  be a time homogeneous, aperiodic and Harris recurrent Markov chain with transition kernel  $P$  and limiting invariant probability distribution  $\pi$ .

Consider the estimator along one walk without burn-in, i.e.

$$\hat{I}_n = \frac{1}{n} \sum_{i=0}^{n-1} g(X_i) \tag{4.1}$$

of the unknown value  $I$ . By Theorem 2.1.5,  $\hat{I}_n \rightarrow I$ , as  $n \rightarrow \infty$ , with probability 1. Moreover, assume for a moment that a  $\sqrt{n}$ -CLT holds and let  $\sigma_g^2$  be the asymptotic variance, as defined in (3.2).

We will study the following sequential procedure. Let  $n^* = n^*(\varepsilon)$  be the first time that

$$q_\bullet \frac{\hat{\sigma}_n}{\sqrt{n}} + p(n) \leq \varepsilon, \quad (4.2)$$

where  $\hat{\sigma}_n^2$  is an estimate of  $\sigma_g^2$  at time  $n$ , and  $q_\bullet$  is an appropriate quantile,  $p(n) > 0$  is a strictly positive decreasing function on  $Z_+$ , and  $\varepsilon > 0$  is the desired half-width.

At time  $n^*$  we build an interval  $I^*(\varepsilon) := [\hat{I}_{n^*} - \varepsilon, \hat{I}_{n^*} + \varepsilon]$  of width  $2\varepsilon$ . For independent samples such procedures are known to work well and belong to classical results of sequential statistics (c.f. [Chow & Robbins 1965], [Nadas 1969] and [Liu, W 1997]). However in our context we have to apply the following result from [Glynn & Whitt 1992].

**Theorem 4.1.1** (Glynn & Whitt 1992). *If*

- (a) *A functional central limit theorem holds, i.e. as  $n \rightarrow \infty$ , the distribution of*

$$Y_n(t) := \frac{1}{\sqrt{n}} \sum_{i=1}^{\lfloor nt \rfloor} g(X_i)$$

*converges to Brownian motion with variance  $\sigma_g^2$  weakly in the Skorohod space on any finite interval,*

- (b)  $\hat{\sigma}_n^2 \rightarrow \sigma_g^2$  *with probability 1 as  $n \rightarrow \infty$ ,*  
(c) *The sequence  $p(n)$  is strictly positive and decreasing and  $p(n) = o(n^{-1/2})$ ,*

*then*

$$P(I \in I^*(\varepsilon)) \rightarrow 1 - \delta, \quad \text{as } \varepsilon \rightarrow 0. \quad (4.3)$$

Markov chains often enjoy a functional central limit theorem under the same conditions that ensure the standard  $\sqrt{n}$ -CLT. In particular the following results are well known:

**Theorem 4.1.2.** *Assume  $(X_n)_{n \geq 0}$  is a Harris ergodic Markov chain. If one of the following conditions holds, then a functional central limit theorem also holds.*

- (a) (due to [Doukhan et al. 1994]) The chain is geometrically ergodic and  $E_\pi[g^2(x)(\log^+ |g(x)|)] < \infty$ ,
- (b) (due to [Roberts & Rosenthal 1997b]) The chain is geometrically ergodic, reversible, and  $E_\pi g^2(x) < \infty$ ,
- (c) (due to [Billingsley 1968]) The chain is uniformly ergodic and  $E_\pi g^2(x) < \infty$ .

The goal of this chapter is to obtain additionally condition (b) of Theorem 4.1.1 for a suitable estimator  $\hat{\sigma}_n^2$  of  $\sigma_g^2$ , under possibly weak assumptions and consequently conclude (4.3). In particular we will need stronger assumptions than those listed in Theorem 4.1.2, thus condition (a) of Theorem 4.1.1 will hold automatically.

## 4.2 Estimating Asymptotic Variance

We will discuss two methods of estimating the asymptotic variance described in [Jones et al. 2006], based on batch means and regenerative simulation.

### 4.2.1 Batch Means

For the bath means estimator suppose that  $n - 1$  iterations of the algorithm are performed and we partition the trajectory of length  $n$  into  $a_n$  blocks of length  $b_n$  i.e.

$$n \simeq a_n b_n$$

Define  $\bar{Y}_1, \dots, \bar{Y}_{a_n}$  as

$$\bar{Y}_j := \frac{1}{b_n} \sum_{i=(j-1)b_n}^{jb_n-1} g(X_i).$$

Then the bath means estimate of  $\sigma_g^2$  is

$$\hat{\sigma}_{BM}^2 = \frac{b_n}{a_n - 1} \sum_{j=1}^{a_n} (\bar{Y}_j - \hat{I}_n)^2. \quad (4.4)$$

In the next section we provide an appropriate strategy for choosing  $a_n$  and  $b_n$  for  $\hat{\sigma}_{BM}^2$  to be a consistent estimator.

## 4.2.2 Regenerative Estimation

Assume that the following minorization condition with  $m = 1$ , as introduced in Definition 2.2.2 holds.

$$P(x, \cdot) \geq s(x)\nu(\cdot), \quad \text{for all } x \in \mathcal{X}, \quad (4.5)$$

and define the residual transition kernel  $R(x, dy)$  as

$$R(x, dy) := \begin{cases} (1 - s(x))^{-1}(P(x, dy) - s(x)\nu(dy)) & \text{if } s(x) < 1, \\ 0 & \text{if } s(x) = 1. \end{cases}$$

By straightforward modification of the split chain construction of Section 2.2 we obtain a bivariate process  $(X_n, Y_n)_{n \geq 0}$  that evolves according to the following transition rule:

- given  $X_n = x$ , draw  $Y_n \sim \text{Bernoulli}(s(x))$
- If  $Y_n = 1$ , then draw  $X_{n+1} \sim \nu(\cdot)$ , otherwise draw  $X_{n+1} \sim R(x, \cdot)$ .

Moreover, the artificial atom  $\check{\alpha}$  is now of the form  $\check{\alpha} = \mathcal{X} \times \{1\}$ . Let us simplify the notation of Section 2.2 by setting  $\tau_n = \tau_{\check{\alpha}}(n)$ , for  $n = 1, 2, \dots$ . Suppose also that  $X_0 \sim \nu$  and set  $\tau_0 = -1$  to keep notation coherent with probabilistic behavior of the chain. Define also  $N_i = \tau_{i+1} - \tau_i$ , for  $i = 0, 1, \dots$ , and recall  $s_i$  defined by (2.19). Since  $m = 1$ ,

$$s_i = \sum_{j=\tau_i+1}^{\tau_{i+1}} g(X_j),$$

and observe that the  $(N_i, s_i)$  pairs are iid random variables.

For regenerative estimation of the asymptotic variance we will need  $(Y_i)_{i \geq 0}$ , thus we must simulate the split chain  $(X_i, Y_i)_{i \geq 0}$ , not only the initial chain  $(X_i)_{i \geq 0}$ . However the simulation from  $R(x, \cdot)$  in real life examples is often challenging. The following solution to this problem is provided in [Mykland et al. 1995].

Suppose that  $P(x, \cdot)$  has a density  $k(\cdot|x)$  and  $\nu(\cdot)$  has a density  $v(\cdot)$  with respect to a reference measure  $dx$ . Given  $X_i = x$  draw  $X_{i+1} \sim k(\cdot|x)$  and draw  $Y_i$  from the distribution of  $Y_i|X_i, X_{i+1}$ , that is

$$Y_i \sim \text{Bernoulli}\left(\frac{s(X_i)v(X_{i+1})}{k(X_{i+1}|X_i)}\right).$$

The method is feasible in many settings of practical interest (cf. [Mykland et al. 1995], [Jones et al. 2006]).

Once we are able to simulate the split chain  $(X_i, Y_i)_{i \geq 0}$ , we can observe  $\tau_0, \tau_1, \dots$  and compute the following regenerative estimator of  $I$ .

$$\hat{I}_{\tau_R} = \frac{1}{\tau_R + 1} \sum_{j=0}^{\tau_R} g(X_j), \quad (4.6)$$

where the fixed number  $R$  is the total number of regenerations observed. Note that  $\hat{I}_{\tau_R}$  is a sum of fixed number of iid. random variables. Thus if  $E_\nu N_0^2 < \infty$  and  $E_\nu s_0^2 < \infty$  then

$$\sqrt{R}(\hat{I}_{\tau_R} - I) \rightarrow N(0, \xi_g^2), \quad \text{as } R \rightarrow \infty, \quad (4.7)$$

where

$$\xi_g^2 = \frac{E_\nu(s_0 - N_0 I)^2}{(E_\nu N_0)^2}.$$

Let  $\bar{N} = R^{-1}(\tau_R + 1) = R^{-1} \sum_{i=0}^{R-1} N_i$ . As an approximation for  $\xi_g^2$  one can take the following regenerative estimator

$$\hat{\xi}_{RS}^2 := \frac{1}{R\bar{N}^2} \sum_{i=0}^{R-1} (s_i - \hat{I}_{\tau_R} N_i)^2. \quad (4.8)$$

Now observe that

$$\begin{aligned} \hat{\xi}_{RS}^2 - \xi_g^2 &= \frac{1}{R\bar{N}^2} \sum_{i=0}^{R-1} (s_i - \hat{I}_{\tau_R} N_i)^2 \pm \frac{E_\nu(s_0 - N_0 I)^2}{\bar{N}^2} - \frac{E_\nu(s_0 - N_0 I)^2}{(E_\nu N_0)^2} \\ &= \frac{1}{R\bar{N}^2} \sum_{i=0}^{R-1} \left[ (s_i - \hat{I}_{\tau_R} N_i)^2 \pm (s_i - N_i I)^2 - E_\nu(s_0 - N_0 I)^2 \right] + \\ &\quad + E_\nu(s_0 - N_0 I)^2 \left[ \frac{1}{\bar{N}^2} - \frac{1}{(E_\nu N_0)^2} \right]. \end{aligned}$$

As noticed in [Jones et al. 2006], repeated application of the strong law of large numbers (with  $R \rightarrow \infty$ ) yields that  $\hat{\xi}_{RS}^2$  is a strongly consistent estimator of  $\xi_g^2$  so it is enough to establish conditions  $E_\nu N_0^2 < \infty$  and  $E_\nu s_0^2 < \infty$  for the fixed width methodology to work. This is deferred to the next section.



Clearly, in this modified regenerative setting an asymptotically valid fixed-width result is obtained by terminating the simulation the first time that

$$q \cdot \frac{\hat{\xi}_{RS}}{\sqrt{R}} + p(R) \leq \varepsilon. \quad (4.9)$$

### 4.3 A Lemma and its Consequences

For geometrically ergodic Markov chains hitting times for sets of positive stationary measure have geometrically decreasing tails. In particular the following lemma is shown in [Hobert et al. 2002].

**Lemma 4.3.1** (Lemma 2 of [Hobert et al. 2002]). *Let  $(X_n)_{n \geq 0}$  be a Harris ergodic chain and assume that (4.5) holds. If  $(X_n)_{n \geq 0}$  is geometrically ergodic, then there exists a  $\beta > 1$ , such that  $E_\pi \beta^{\tau_1} < \infty$ .*

Which immediately yields the following corollary.

**Corollary 4.3.2.** *Under the conditions of Lemma 4.3.1, for any  $a > 0$ ,*

$$\sum_{i=0}^{\infty} \left( P_\pi(\tau_1 \geq i+1) \right)^a \leq \left( E_\pi \beta^{\tau_1} \right)^a \sum_{i=0}^{\infty} \beta^{-a(i+1)} < \infty. \quad (4.10)$$

*Proof.*

$$\begin{aligned} \sum_{i=0}^{\infty} \left( P_\pi(\tau_1 \geq i+1) \right)^a &\leq \sum_{i=0}^{\infty} \left( E_\pi(\mathbb{I}_{\{\tau_1 \geq i+1\}} \beta^{\tau_1} \beta^{-(i+1)}) \right)^a \\ &= \sum_{i=0}^{\infty} \beta^{-a(i+1)} \left( E_\pi(\mathbb{I}_{\{\tau_1 \geq i+1\}} \beta^{\tau_1}) \right)^a \\ &\leq \sum_{i=0}^{\infty} \beta^{-a(i+1)} \left( E_\pi \beta^{\tau_1} \right)^a. \end{aligned}$$

□

Observe also that we can integrate (4.5) with respect to  $\pi$  and obtain  $\pi(\cdot) \geq c\nu(\cdot)$ , where  $c = E_\pi s$ . Thus for any function  $h : \mathcal{X}^\infty \rightarrow R$ ,

$$E_\pi |h(X_0, X_1, \dots)| \geq c E_\nu |h(X_0, X_1, \dots)|. \quad (4.11)$$

Now we are in a position to prove our key result, namely the following lemma.

**Lemma 4.3.3.** *Let  $(X_n)_{n \geq 0}$  be a Harris ergodic Markov chain, assume the minorization condition (4.5) holds, and  $(X_n)_{n \geq 0}$  is geometrically ergodic. Let  $g : \mathcal{X} \rightarrow \mathbb{R}$  be a real valued Borel function. Then, if*

$$E_\pi |g|^{p+\delta} < \infty \quad \text{for some } p > 0 \quad \text{and} \quad \delta > 0,$$

then

$$E_\nu N_0^p < \infty \quad \text{and} \quad E_\nu |s_0|^p < \infty.$$

*Remark 4.3.4.* Lemma 4.3.3 improves the two following results:

- Theorem 2 of [Hobert et al. 2002] that provides the implication

$$E_\pi |g|^{2+\delta} < \infty \Rightarrow E_\nu N_0^2 < \infty \quad \text{and} \quad E_\nu |s_0|^2 < \infty.$$

- Lemma 1 of [Jones et al. 2006] that for  $p \geq 1$  provides implications

$$E_\pi |g|^{2(p-1)+\delta} < \infty \Rightarrow E_\nu N_0^p < \infty \quad \text{and} \quad E_\nu |s_0|^p < \infty.$$

and

$$E_\pi |g|^{2p+\delta} < \infty \Rightarrow E_\nu N_0^p < \infty \quad \text{and} \quad E_\nu |s_0|^{p+\delta} < \infty.$$

*Remark 4.3.5.* Without additional restrictions  $E_\pi |g|^p < \infty$  does not imply  $E_\nu |s_0|^p < \infty$ , so Lemma 4.3.3 can not be improved. To see this note that Theorem 3.3.1 of Chapter 3 combined with the presumption that in the setting of Lemma 4.3.3  $E_\pi |g|^p < \infty$  implies  $E_\nu |s_0|^p < \infty$  yields the Central Limit Theorem for normalized sums of  $g(X_i)$  for geometrically ergodic Markov chains assuming only  $E_\pi g^2 < \infty$ . This however is not enough for the CLT, Bradley in [Bradley 1983] and also Häggström in [Häggström 2005] provide counterexamples. Hence to obtain the implication  $E_\pi |g|^p < \infty \Rightarrow E_\nu |s_0|^p < \infty$ , one needs stronger assumptions, e.g. if  $p = 2$  then uniform ergodicity is enough, as proved in Chapter 3.

*Proof of Lemma 4.3.3.* First note that by (4.11) it is enough to show that

$$E_\pi N_0^p < \infty \quad \text{and} \quad E_\pi |s_0|^p < \infty.$$

Moreover, since  $\max_k \left\{ \frac{k^p}{\beta^k} \right\} < \infty$  for every  $p > 0$  and  $\beta > 1$ , by Lemma 4.3.1 we obtain immediately  $E_\pi N_0^p < \infty$ . Thus we proceed to show that  $E_\pi |s_0|^p < \infty$ . To this end first note that

$$C := \left( (E_\pi |g(X_i)|^{p+\delta})^{\frac{p}{p+\delta}} \right)^{1/p} < \infty. \quad (4.12)$$

For  $p \geq 1$  we use first the triangle inequality in  $L^p$ , then Hölder inequality, then (4.12) and finally Corollary 4.3.2.

$$\begin{aligned}
(E_\pi |s_0|^p)^{1/p} &\leq \left[ E_\pi \left( \sum_{i=0}^{\tau_1} |g(X_i)| \right)^p \right]^{1/p} \\
&= \left[ E_\pi \left( \sum_{i=0}^{\infty} \mathbf{1}(i \leq \tau_1) |g(X_i)| \right)^p \right]^{1/p} \\
&\leq \sum_{i=0}^{\infty} \left[ E_\pi \mathbf{1}(i \leq \tau_1) |g(X_i)|^p \right]^{1/p} \\
&\leq \sum_{i=0}^{\infty} \left[ (E_\pi \mathbf{1}(i \leq \tau_1))^{\frac{\delta}{p+\delta}} (E_\pi |g(X_i)|^{p+\delta})^{\frac{p}{p+\delta}} \right]^{1/p} \\
&= C \sum_{i=0}^{\infty} (P_\pi(\tau_1 \geq i))^{\frac{\delta}{p+\delta}} < \infty. \tag{4.13}
\end{aligned}$$

For  $0 < p < 1$  we use the fact  $x^p$  is concave and then proceed similarly as in (4.13) to obtain

$$\begin{aligned}
E_\pi |s_0|^p &\leq E_\pi \left( \sum_{i=0}^{\infty} \mathbf{1}(i \leq \tau_1) |g(X_i)| \right)^p \\
&\leq \sum_{i=0}^{\infty} E_\pi \mathbf{1}(i \leq \tau_1) |g(X_i)|^p \\
&\leq C^p \sum_{i=0}^{\infty} (P_\pi(\tau_1 \geq i))^{\frac{\delta}{p+\delta}} < \infty.
\end{aligned}$$

□

Lemma 4.3.3 allows us to restate results from section 3.2 of [Jones et al. 2006] with relaxed assumptions. In particular in Lemma 2 and in Proposition 3 therein it is enough to assume  $E_\pi |g|^{2+\delta+\varepsilon} < \infty$  for some  $\delta > 0$  and some  $\varepsilon > 0$ , instead of  $E_\pi |g|^{4+\delta} < \infty$  for some  $\delta > 0$ . Modifications of the (rather long and complicated) proofs in [Jones et al. 2006] are straightforward. Hence we have

**Lemma 4.3.6** (Part b of Lemma 2 of [Jones et al. 2006]). *Let  $(X_n)_{n \geq 0}$  be a Harris ergodic Markov chain with invariant distribution  $\pi$ . If  $(X_n)_{n \geq 0}$  is*

geometrically ergodic, (4.5) holds and  $E_\pi |g|^{2+\delta+\varepsilon} < \infty$  for some  $\delta > 0$  and some  $\varepsilon > 0$ , then there exists a constant  $0 < \sigma_g < \infty$ , and a sufficiently large probability space such that

$$\left| \sum_{i=1}^n g(X_i) - nE_\pi g - \sigma_g B(n) \right| = O(\gamma(n))$$

with probability 1 as  $n \rightarrow \infty$ , where  $\gamma(n) = n^\alpha \log n$ ,  $\alpha = 1/(2 + \delta)$ , and  $B = \{B(t), t \geq 0\}$  denotes a standard Brownian motion.

**Proposition 4.3.7** (Proposition 3 of [Jones et al. 2006]). *Let  $(X_n)_{n \geq 0}$  be a Harris ergodic Markov chain with invariant distribution  $\pi$ . Further, suppose  $(X_n)_{n \geq 0}$  is geometrically ergodic, (4.5) holds and  $E_\pi |g|^{2+\delta+\varepsilon} < \infty$  for some  $\delta > 0$  and some  $\varepsilon > 0$ . If*

1.  $a_n \rightarrow \infty$ , as  $n \rightarrow \infty$ ,
2.  $b_n \rightarrow \infty$  and  $b_n/n \rightarrow 0$  as  $n \rightarrow \infty$ ,
3.  $b_n^{-1} n^{2\alpha} [\log n]^3 \rightarrow 0$  as  $n \rightarrow \infty$ , where  $\alpha = 1/(2 + \delta)$ ,
4. there exists a constant  $c \geq 1$ , such that  $\sum_{n=1}^{\infty} (b_n/n)^c < \infty$ ,

Then  $\hat{\sigma}_{BM}^2 \rightarrow \sigma_g^2$  w.p.1 as  $n \rightarrow \infty$ .

*Concluding Remark 4.3.8.* Compare the foregoing result with Section 4.2.2 or with Proposition 1 of [Jones et al. 2006] to see that both methods described here, i.e. regenerative simulation (RS) and batch means (CBM), provide strongly consistent estimators of  $\sigma_g^2$  under the same assumption for the target function  $g$ .

# Chapter 5

## Fixed-Width Nonasymptotic Results under Drift Condition

In this Chapter we establish nonasymptotic fixed width estimation. We assume a drift condition towards a small set and bound the mean square error of estimators obtained by taking averages along a single trajectory of a Markov chain Monte Carlo algorithm. We use these bounds to determine the length of the trajectory and the burn-in time that ensures  $(\varepsilon - \alpha)$ -approximation, i.e. desired precision of estimation with given probability. Let  $I$  be the value of interest and  $\hat{I}$  its MCMC estimate. Precisely, our lower bounds for the length of the trajectory and burn-in time ensure that

$$P(|\hat{I} - I| \leq \varepsilon) \geq 1 - \alpha$$

and depend only and explicitly on drift parameters,  $\varepsilon$  and  $\alpha$ . Next we introduce an MCMC estimator based on the median of multiple shorter runs. It turns out that this estimation scheme allows for sharper bounds for the total simulation cost required for the  $(\varepsilon - \alpha)$ -approximation. For both estimation schemes numerical examples are provided that include practically relevant Gibbs samplers for a hierarchical random effects model.

### 5.1 Introduction

Recall the estimation strategies introduced in Section 1.2 and described by (1.3-1.5). *Estimation Along one Walk* uses average along a single trajectory of the underlying Markov chain and discards the initial part to reduce bias.

The estimate of the unknown value  $I = \int_{\mathcal{X}} f(x)\pi(dx)$  is of the form

$$\hat{I}_{t,n} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i) \quad (5.1)$$

and  $t$  is called the burn-in time.

The strategy is believed to be more efficient than *estimation along one walk with spacing* and *multiple run* described in Section 1.2 and is usually the practitioners choice. Some precise results are available for reversible Markov chains. Geyer in [Geyer 1992] shows that using spacing as in (1.4) is ineffective (in terms of asymptotic variance) and Chan and Yue in [Chan & Yue 1996] prove that (5.1) is asymptotically efficient in a class of linear estimators (in terms of mean square error).

The goal of this chapter is to derive lower bounds for  $n$  and  $t$  in (5.1), that minimize the total computation cost  $n+t$ , and that ensure the following condition of  $(\varepsilon, \alpha)$ -approximation:

$$P(|\hat{I}_{t,n} - I| \leq \varepsilon) \geq 1 - \alpha, \quad (5.2)$$

where  $\varepsilon$  is the precision of estimation and  $1 - \alpha$ , the confidence level. Due to results in [Geyer 1992] and [Chan & Yue 1996] no other linear modifications of the estimation scheme in (5.1) are analyzed. To decrease the total simulation cost for (5.2) we introduce instead a nonlinear estimator based on the median of multiple shorter runs.

Results of this or related type have been obtained for discrete state space  $\mathcal{X}$  and bounded target function  $f$  by Aldous in [Aldous 1987], Gillman in [Gillman 1998] and recently by León and Perron in [León & Perron 2004]. Niemiro and Pokarowski in [Niemiro & Pokarowski 2007] give results for relative precision estimation. For uniformly ergodic chains on continuous state space  $\mathcal{X}$  and bounded function  $f$ , Hoeffding type inequalities are available (due to Glynn and Ormonait in [Glynn & Ormoneit 2002], and an improved bound due to Meyn et al. in [Kontoyiannis et al. 2005]) and can easily lead to the desired  $(\varepsilon, \alpha)$ -approximation. To our best knowledge there are no explicit bounds for  $n$  and  $t$  in more general settings, especially when  $f$  is not bounded and the chain is not uniformly ergodic. A remarkable presentation of the state of the art approach to dealing with this problem is provided by Jones et al. in the recent paper [Jones et al. 2006]. They suggest two procedures for constructing consistent estimators for the variance of the asymptotic

normal distribution for geometrically ergodic split chains and thus under the additional assumption of  $E_\pi |f|^{2+\delta} < \infty$  for some  $\delta > 0$  (see Chapter 4 here for this weakened assumption and details of the procedure).

Our approach is to assume a version of the well known drift condition towards a small set (Assumption 5.2.1) and give explicit lower bounds on  $n$  and  $t$  in terms of drift parameters defined in Assumption 5.2.1 and approximation parameters defined in (5.2).

The rest of the Chapter is organized as follows. In Section 5.2 we introduce the drift condition assumption and preliminary results. In Section 5.3 we obtain an explicit bound for the mean square error of the estimator defined in (1.3). In Section 5.4 we construct two different  $(\varepsilon - \alpha)$ -approximation procedures, one based on the sample mean of one long trajectory and the other based on the median of multiple shorter runs. We close with examples in Sections 5.5 and 5.6, in particular we show how to obtain explicit lower bounds for  $t$  and  $n$  that guarantee the  $\varepsilon - \alpha$ -approximation for a hierarchical random effects model of practical relevance.

## 5.2 A Drift Condition and Preliminary Lemmas

Since in what follows we deal with integrals of unbounded functions  $f$  with respect to probability measures, the very common total variation distance defined by (2.1) is inappropriate for measuring distances between probability measures and we need to use the  $V$ -norm and  $V$ -norm distance introduced in Section 2.1.

We analyze the MCMC estimation along a single trajectory under the following assumption of a drift condition towards a small set.

**Assumption 5.2.1.** (A.1) *Small set.* There exist  $C \in \mathcal{B}(\mathcal{X})$ ,  $\tilde{\beta} > 0$  and a probability measure  $\nu$  on  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  such that for all  $x \in C$  and  $A \in \mathcal{B}(\mathcal{X})$

$$P(x, A) \geq \tilde{\beta}\nu(A).$$

(A.2) *Drift.* There exist a function  $V : \mathcal{X} \rightarrow [1, \infty)$  and constants  $\lambda < 1$  and  $K < \infty$  satisfying

$$PV(x) \leq \begin{cases} \lambda V(x), & \text{if } x \notin C, \\ K, & \text{if } x \in C. \end{cases}$$

(A.3) *Aperiodicity.* There exists  $\beta > 0$  such that  $\tilde{\beta}\nu(C) \geq \beta$ .

In the sequel we refer to  $\tilde{\beta}, V(x), \lambda, K, \beta$  as drift parameters.

*Remark 5.2.2.* Establishing a drift condition for real life examples is usually not an easy task. As indicated in [Meyn & Tweedie 1993] polynomials are often suitable candidates for a drift function  $V$  and also functions proportional to  $\pi^{1/2}$  may turn out to be a lucky choice. Computable toy and real life examples of [Baxendale 2005] and [Jones & Hobert 2004] confirm this observations.

*Remark 5.2.3.* There is a strong probabilistic intuition behind Assumption 5.2.1. Every time the chain visits the small set  $C$ , it regenerates with probability  $\tilde{\beta}$ . The role of the drift condition (A.2) is to guarantee that the chain visits the small set  $C$  frequently enough. Typically  $C$  is in the „center” of the state space  $\mathcal{X}$  and the drift function  $V$  takes small values on  $C$  and increases as it goes away from  $C$ . Assume first that  $X_n = x \notin C$ . The condition  $PV(x) \leq \lambda V(x)$  means that  $X_{n+1} \sim P(x, \cdot)$  is on average getting closer to  $C$  (closer in terms of  $V$ ). Whereas  $PV(x) \leq K$  for  $X_n = x \in C$  means that  $X_{n+1}$  will perhaps jump out of  $C$ , but not too far away, i.e. the integral of  $V$  with respect to the distribution of  $X_{n+1}$  is bounded (by the same value) for all  $x \in C$ . Assumption (A.3) together with (A.1) imply aperiodicity.

Assumption 5.2.1 is often used and widely discussed in Markov chains literature. Substantial effort has been devoted to establishing convergence rates for Markov chains under the drift condition (A.1-3) or related assumptions. For discussion of various drift conditions and their relation see Meyn and Tweedie [Meyn & Tweedie 1993]. For quantitative bounds on convergence rates of Markov chains see the survey paper by Roberts and Rosenthal [Roberts & Rosenthal 2005] and references therein. In the sequel we make use of the recent convergence bounds obtained by Baxendale in [Baxendale 2005].

**Theorem 5.2.4** (Baxendale [Baxendale 2005]). *Under Assumption 5.2.1  $(X)_{n \geq 0}$  has a unique stationary distribution  $\pi$  and  $\pi V < \infty$ . Moreover, there exists  $\rho < 1$  depending only and explicitly on  $\tilde{\beta}, \beta, \lambda$  and  $K$  such that whenever  $\rho < \gamma < 1$  there exists  $M < \infty$  depending only and explicitly on  $\gamma, \tilde{\beta}, \beta, \lambda$  and  $K$  such that for all  $n \geq 0$*

$$\| \| P^n - \pi \| \|_V \leq M \gamma^n. \tag{5.3}$$



When we refer in the sequel to  $V$ -uniform ergodicity, we mean the convergence determined by (5.3). There are different formulas for  $\rho$  and  $M$  for general operators, self adjoint operators and self adjoint positive operators in both atomic and nonatomic case. We give them in Section 5.8 for the sake of completeness. To our knowledge the above-mentioned theorem gives the best available explicit constants.

**Corollary 5.2.5.** *Under Assumption 5.2.1*

$$\|\pi_0 P^n - \pi\|_V \leq \min\{\pi_0 V, \|\pi_0 - \pi\|_V\} M \gamma^n,$$

where  $M$  and  $\gamma$  are such as in Theorem 5.2.4.

*Proof.* From Theorem 5.2.4 we have  $\|P^n(x, \cdot) - \pi(\cdot)\|_V \leq M \gamma^n V(x)$ , which yields

$$\begin{aligned} \pi_0 V M \gamma^n &\geq \int_{\mathcal{X}} \|P^n(x, \cdot) - \pi(\cdot)\|_V \pi_0(dx) \geq \sup_{|g| \leq V} \int_{\mathcal{X}} |P^n(x, \cdot)g - \pi g| \pi_0(dx) \\ &\geq \sup_{|g| \leq V} |\pi_0 P^n g - \pi g| = \|\pi_0 P^n - \pi\|_V. \end{aligned}$$

Now let  $b_V = \inf_{x \in \mathcal{X}} V(x)$ . Since  $\|\cdot\|_V$  is an operator norm and  $\pi$  is invariant for  $P$ , we have

$$\begin{aligned} \|\pi_0 P^n - \pi\|_V &= b_V \|\pi_0 P^n - \pi\|_V = b_V \|(\pi_0 - \pi)(P^n - \pi)\|_V \\ &\leq b_V \|\pi_0 - \pi\|_V \|P^n - \pi\|_V = \|\pi_0 - \pi\|_V \|P^n - \pi\|_V. \\ &\leq \|\pi_0 - \pi\|_V M \gamma^n. \end{aligned}$$

□

Now we focus on the following simple but useful observation.

**Lemma 5.2.6.** *If for a Markov chain  $(X_n)_{n \geq 0}$  on  $\mathcal{X}$  with transition kernel  $P$  Assumption 5.2.1 holds with parameters  $\tilde{\beta}, V(x), \lambda, K, \beta$ , it holds also with  $\tilde{\beta}_r := \tilde{\beta}, V_r(x) := V(x)^{1/r}, \lambda_r := \lambda^{1/r}, K_r := K^{1/r}, \beta_r := \beta$  for every  $r > 1$ .*

*Proof.* It is enough to check (A.2). For  $x \notin C$  by Jensen inequality we have

$$\lambda V(x) \geq \int_{\mathcal{X}} V(y) P(x, dy) \geq \left( \int_{\mathcal{X}} V(y)^{1/r} P(x, dy) \right)^r$$

and hence  $PV(x)^{1/r} \leq \lambda^{1/r} V(x)^{1/r}$ , as claimed. Similarly for  $x \in C$  we obtain  $PV(x)^{1/r} \leq K^{1/r}$ . □

Lemma 5.2.6 together with Theorem 5.2.4 yield the following corollary.

**Corollary 5.2.7.** *Under Assumption 5.2.1 we have*

$$\|P^n - \pi\|_{V^{1/r}} \leq M_r \gamma_r^n,$$

where  $M_r$  and  $\gamma_r$  are constants defined as in Theorem 5.2.4 resulting from drift parameters defined in Lemma 5.2.6.

Integrating the drift condition with respect to  $\pi$  yields the following bound on  $\pi V$ .

**Lemma 5.2.8.** *Under Assumption 5.2.1*

$$\pi V \leq \pi(C) \frac{K - \lambda}{1 - \lambda} \leq \frac{K - \lambda}{1 - \lambda}.$$

Let  $f_c = f - \pi f$ . The next lemma provides a bound on  $\|f_c\|^p|_V$  in terms of  $\|f\|^p|_V$  without additional effort.

**Lemma 5.2.9.** *Under Assumption 5.2.1*

$$\|f_c\|^p|_V \leq \left( C_{f_V}^{1/p} + \frac{\pi(C)}{b_V^{1/p}} K_{p,\lambda} \right)^2 \leq \left( C_{f_V}^{1/p} + K_{p,\lambda} \right)^2,$$

where  $b_V = \inf_{x \in \mathcal{X}} V(x)$ ,  $C_{f_V} = \|f\|^p|_V$  and  $K_{p,\lambda} = \frac{K^{1/p} - \lambda^{1/p}}{1 - \lambda^{1/p}}$ .

*Proof.* Note that  $\pi V^{1/p} \leq \pi(C) K_{p,\lambda} \leq K_{p,\lambda}$  by Lemma 5.2.8 and proceed:

$$\begin{aligned} \|f_c\|^p|_V &= \sup_{x \in \mathcal{X}} \frac{|f(x) - \pi f|^p}{V(x)} \leq \sup_{x \in \mathcal{X}} \frac{\left( C_{f_V}^{1/p} V^{1/p}(x) + \pi|f| \right)^p}{V(x)} \\ &\leq \sup_{x \in \mathcal{X}} \frac{\left( C_{f_V}^{1/p} V^{1/p}(x) + \pi(C) K_{p,\lambda} \right)^p}{V(x)} \leq C_{f_V}^p \left( 1 + \frac{\pi(C) K_{p,\lambda}}{b_V^{1/p} C_{f_V}^{1/p}} \right)^p. \end{aligned}$$

□

## 5.3 MSE Bounds

By  $MSE(\hat{I}_{0,n})$  we denote the mean square error of  $\hat{I}_{0,n}$ , i.e.

$$MSE(\hat{I}_{0,n}) = E_{\pi_0}[\hat{I}_{0,n} - I]^2.$$

Bonds on  $MSE(\hat{I}_{0,n})$  are essential to establish  $(\varepsilon - \alpha)$ -approximation of type (5.2) and are also of independent interest.

**Theorem 5.3.1.** *Assume the Drift Condition 5.2.1 holds and  $X_0 \sim \pi_0$ . Then for every measurable function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , every  $p \geq 2$  and every  $r \in [\frac{p}{p-1}, p]$*

$$MSE(\hat{I}_{0,n}) \leq \frac{\|f_c\|^p|_V^{2/p}}{n} \left(1 + \frac{2M_r\gamma_r}{1 - \gamma_r}\right) \left(\pi V + \frac{M \min\{\pi_0 V, \|\pi_0 - \pi\|_V\}}{n(1 - \gamma)}\right), \quad (5.4)$$

where  $f_c = f - \pi f$  and constants  $M, \gamma, M_r, \gamma_r$  depend only and explicitly on  $\tilde{\beta}, \beta, \lambda$  and  $K$  from Assumption 5.2.1 as in Theorem 5.2.4 and Corollary 5.2.6.

The formulation of the foregoing Theorem 5.3.1 is motivated by a trade-off between small  $V$  and small  $\lambda$  in Assumption 5.2.1. It should be intuitively clear that establishing the drift condition for a quickly increasing  $V$  should result in smaller  $\lambda$  at the cost of bigger  $\pi V$ . So it may be reasonable to look for a valid drift condition with  $V \geq C\|f_c\|^p$  for some  $p > 2$  instead of the natural choice of  $p = 2$ . Lemma 5.2.6 should strengthen this intuition. The most important special case for  $p = r = 2$  is emphasized below as a corollary.

The unknown value  $\pi_0 V$  in (5.4) depends on  $\pi_0$  which is users choice and usually a deterministic point. Also, in many cases a fairly small bound for  $\pi V$  should be possible to obtain by direct calculations, since in the typical setting  $\pi$  is exponentially concentrated whereas  $V$  is a polynomial of degree 2. These calculations should probably borrow from those used to obtain the minorization and drift conditions. However, in absence of a better bound for  $\pi V$  Lemma 5.2.8 is at hand. Similarly Lemma 5.2.9 bounds the unknown value  $\|f_c\|^p|_V^{2/p}$  in terms of  $\|f\|^p|_V$ . Note that in applications both  $f$  and  $V$  have explicit formulas known to the user and  $\|f\|^p|_V$  can be evaluated directly or easily bounded.

*Proof.* Note that  $\|f\|_{V^{1/r}}^r = \|f\|^r|_V$ . Without loss of generality consider  $f_c$  instead of  $f$  and assume  $\|f_c\|^p|_V = 1$ . In this setting  $\|f_c\|^2|_V \leq 1$ ,  $Var_{\pi} f_c =$

$\pi f_c^2 \leq \pi V$ ,  $MSE(\hat{I}_{0,n}) = E_{\pi_0}(\hat{I}_{0,n})^2$ , and also for every  $r \in [\frac{p}{p-1}, p]$ ,

$$|f_c|_{V^{1/r}} \leq \|f_c\|_{V^{1/r}}^{p/r} = 1 \quad \text{and} \quad |f_c|_{V^{1-1/r}} \leq \|f_c\|_{V^{1-1/r}}^{p-p/r} = 1.$$

Obviously

$$nMSE(\hat{I}_{0,n}) = \frac{1}{n} \sum_{i=0}^{n-1} E_{\pi_0} f_c(X_i)^2 + \frac{2}{n} \sum_{i=0}^{n-2} \sum_{j=i+1}^{n-1} E_{\pi_0} f_c(X_i) f_c(X_j). \quad (5.5)$$

We start with a bound for the first term of the right hand side of (5.5). Since  $f_c^2(x) \leq V(x)$ , we use Corollary 5.2.5 for  $f_c^2$ . Let  $C = \min\{\pi_0 V, \|\pi_0 - \pi\|_V\}$  and proceed

$$\frac{1}{n} \sum_{i=0}^{n-1} E_{\pi_0} f_c(X_i)^2 = \frac{1}{n} \sum_{i=0}^{n-1} \pi_0 P^i f_c^2 \leq \pi f_c^2 + \frac{1}{n} \sum_{i=0}^{n-1} CM \gamma^i \leq \pi V + \frac{CM}{n(1-\gamma)}. \quad (5.6)$$

To bound the second term of the right hand side of (5.5) note that  $|f_c| \leq V^{1/r}$  and use Corollary 5.2.7.

$$\begin{aligned} \frac{2}{n} \sum_{i=0}^{n-2} \sum_{j=i+1}^{n-1} E_{\pi_0} f_c(X_i) f_c(X_j) &= \frac{2}{n} \sum_{i=0}^{n-2} \sum_{j=i+1}^{n-1} \pi_0 (P^i (f_c P^{j-i} f_c)) \\ &\leq \frac{2}{n} \sum_{i=0}^{n-2} \sum_{j=i+1}^{n-1} \pi_0 (P^i (|f_c| P^{j-i} |f_c|)) \\ &\leq \frac{2M_r}{n} \sum_{i=0}^{n-2} \sum_{j=i+1}^{\infty} \gamma_r^{j-i} \pi_0 (P^i (|f_c| V^{1/r})) \\ &\leq \frac{2M_r \gamma_r}{n(1-\gamma_r)} \sum_{i=0}^{n-2} \pi_0 (P^i (|f_c| V^{1/r})) = \spadesuit \end{aligned}$$

Since  $|f_c| \leq V^{1/r}$  and  $|f_c| \leq V^{1-1/r}$ , also  $|f_c V^{1/r}| \leq V$  and we use Corollary 5.2.5 for  $|f_c| V^{1/r}$ .

$$\spadesuit \leq \frac{2M_r \gamma_r}{n(1-\gamma_r)} \sum_{i=0}^{n-2} (\pi (|f_c| V^{1/r}) + CM \gamma^i) \leq \frac{2M_r \gamma_r}{1-\gamma_r} \left( \pi V + \frac{CM}{n(1-\gamma)} \right). \quad (5.7)$$

Combine (5.6) and (5.7) to obtain

$$MSE(\hat{I}_{0,n}) \leq \frac{\|f_c\|_V^{2/p}}{n} \left( 1 + \frac{2M_r \gamma_r}{1-\gamma_r} \right) \left( \pi V + \frac{CM}{n(1-\gamma)} \right).$$

□

**Corollary 5.3.2.** *In the setting of Theorem 5.3.1, we have in particular*

$$MSE(\hat{I}_{0,n}) \leq \frac{|f_c^2|_V}{n} \left( 1 + \frac{2M_2\gamma_2}{1-\gamma_2} \right) \left( \pi V + \frac{M \min\{\pi_0 V, \|\pi_0 - \pi\|_V\}}{n(1-\gamma)} \right). \quad (5.8)$$

The foregoing bound is easy to interpret:  $\pi V |f_c^2|_V$  should be close to  $Var_{\pi} f$  for an appropriate choice of  $V$ , moreover  $2M_2\gamma_2/(1-\gamma_2)$  corresponds to the autocorrelation of the chain and the last term  $M \min\{\pi_0 V, \|\pi_0 - \pi\|_V\}/n(1-\gamma)$  is the price for nonstationarity of the initial distribution. See also Theorem 5.3.4 for further interpretation.

Theorem 5.3.1 is explicitly stated for  $\hat{I}_{0,n}$ , but the structure of the bound is flexible enough to cover most typical settings as indicated below.

**Corollary 5.3.3.** *In the setting of Theorem 5.3.1,*

$$MSE(\hat{I}_{0,n}) \leq \frac{\pi V \|f_c\|_V^{2/p}}{n} \left( 1 + \frac{2M_r\gamma_r}{1-\gamma_r} \right), \quad \text{if } \pi_0 = \pi, \quad (5.9)$$

$$MSE(\hat{I}_{0,n}) \leq \frac{\|f_c\|_V^{2/p}}{n} \left( 1 + \frac{2M_r\gamma_r}{1-\gamma_r} \right) \left( \pi V + \frac{MV(x)}{n(1-\gamma)} \right), \quad \text{if } \pi_0 = \delta_x, \quad (5.10)$$

$$MSE(\hat{I}_{t,n}) \leq \frac{\|f_c\|_V^{2/p}}{n} \left( 1 + \frac{2M_r\gamma_r}{1-\gamma_r} \right) \left( \pi V + \frac{M^2\gamma^t V(x)}{n(1-\gamma)} \right), \quad \text{if } \pi_0 = \delta_x. \quad (5.11)$$

*Proof.* Only (5.11) needs a proof. Note that  $X_t \sim \delta_x P^t$ . Now use Theorem 5.2.4 to see that  $\|\delta_x P^t - \pi\|_V \leq M\gamma^t V(x)$ , and apply Theorem 5.3.1 with  $\pi_0 = \delta_x P^t$ . □

Bound (5.9) corresponds to the situation when a perfect sampler is available. For deterministic start without burn-in and with burn-in (5.10) and (5.11) should be applied respectively.

Next we derive computable bounds for the asymptotic variance  $\sigma_f^2$  in central limit theorems for Markov chains under the assumption of the Drift Condition 5.2.1.

**Theorem 5.3.4.** *Under the Drift Condition 5.2.1 the Markov chain  $(X_n)_{n \geq 0}$  and a function  $f$ , such that  $|f_c^2|_V < \infty$  (or equivalently  $|f^2|_V < \infty$ ), admit a central limit theorem, i.e.:*

$$\sqrt{n}(\hat{I}_{0,n} - I) \xrightarrow{d} N(0, \sigma_f^2) \quad \text{as } n \rightarrow \infty, \quad (5.12)$$

moreover

$$\sigma_f^2 = \lim_{n \rightarrow \infty} n E_\pi [\hat{I}_{0,n} - I]^2 \leq \pi V \|f_c\|_V^{2/p} \left( 1 + \frac{2M_r \gamma_r}{1 - \gamma_r} \right). \quad (5.13)$$

*Proof.* The CLT (i.e. (5.12) and the equation in (5.13)) is a well known fact and results from  $V$ -uniform ergodicity implied by Theorem 5.2.4 combined with Theorems 17.5.4 and 17.5.3 of [Meyn & Tweedie 1993]. Theorem 5.3.1 with  $\pi_0 = \pi$  yields the bound for  $\sigma_f^2$  in (5.13).  $\square$

*Remark 5.3.5.* For reversible Markov chains significantly sharper bounds for  $\sigma_f^2$  can be obtained via functional analytic approach. For a reversible Markov chain its transition kernel  $P$  is a self-adjoint operator on  $L_\pi^2$ . Let  $f \in L_\pi^2$  and  $\pi f = 0$ . If we denote by  $E_f$  the positive measure on  $(-1, 1)$  associated with  $f$  in the spectral decomposition of  $P$ , we obtain (cf. [Kipnis & Varadhan 1986], [Geyer 1992])

$$\sigma_f^2 = \int_{(-1,1)} \frac{1 + \lambda}{1 - \lambda} E_f(d\lambda) \leq \frac{1 + \rho}{1 - \rho} \text{Var}_\pi f \leq \frac{1 + \rho}{1 - \rho} \pi V |f_c^2|_V. \quad (5.14)$$

Where the first inequality in (5.14) holds if we are able to bound the spectral radius of  $P$  acting on  $L_\pi^2$  by some  $\rho < 1$  (cf. [Geyer 1992], [Roberts & Rosenthal 1997b]). Corollary 6.1 of [Baxendale 2005] yields the required bound with  $\rho$  defined as in Theorem 5.2.4.

## 5.4 $(\varepsilon - \alpha)$ -Approximation

$(\varepsilon - \alpha)$ -approximation is an easy corollary of  $MSE$  bounds by the Chebyshev inequality.

**Theorem 5.4.1** ( $(\varepsilon - \alpha)$ -approximation). *Let*

$$b = \frac{\pi V \|f_c\|^p |V|^{2/p}}{\varepsilon^2 \alpha} \left(1 + \frac{2M_r \gamma_r}{1 - \gamma_r}\right), \quad (5.15)$$

$$c = \frac{M \min\{\pi_0 V, \|\pi_0 - \pi\|_V\} \|f_c\|^p |V|^{2/p}}{\varepsilon^2 \alpha (1 - \gamma)} \left(1 + \frac{2M_r \gamma_r}{1 - \gamma_r}\right), \quad (5.16)$$

$$n(t) = \frac{b + \sqrt{b^2 + 4c(t)}}{2}, \quad (5.17)$$

$$c(t) = \frac{M^2 \gamma^t V(x) \|f_c\|^p |V|^{2/p}}{\varepsilon^2 \alpha (1 - \gamma)} \left(1 + \frac{2M_r \gamma_r}{1 - \gamma_r}\right), \quad (5.18)$$

$$\tilde{c} = \frac{M^2 V(x) \|f_c\|^p |V|^{2/p}}{\varepsilon^2 \alpha (1 - \gamma)} \left(1 + \frac{2M_r \gamma_r}{1 - \gamma_r}\right). \quad (5.19)$$

Then under Assumption 5.2.1,

$$P(|\hat{I}_{0,n} - I| \leq \varepsilon) \geq 1 - \alpha, \quad \text{if } X_0 \sim \pi_0, \quad n \geq \frac{b + \sqrt{b^2 + 4c}}{2}. \quad (5.20)$$

$$P(|\hat{I}_{t,n} - I| \leq \varepsilon) \geq 1 - \alpha, \quad \text{if } \begin{cases} X_0 \sim \delta_x, \\ t \geq \max\left\{0, \log_\gamma \left(\frac{2 + \sqrt{4 + b^2 \ln^2 \gamma}}{\tilde{c} \ln^2 \gamma}\right)\right\} \\ n \geq n(t). \end{cases} \quad (5.21)$$

And the above bounds in (5.21) give the minimal length of the trajectory  $(t + n)$  resulting from (5.11).

*Proof.* From the Chebyshev's inequality we get

$$\begin{aligned} P(|\hat{I}_{t,n} - I| \leq \varepsilon) &= 1 - P(|\hat{I}_{t,n} - I| \geq \varepsilon) \\ &\geq 1 - \frac{MSE(\hat{I}_{t,n})}{\varepsilon^2} \geq 1 - \alpha \quad \text{if } MSE(\hat{I}_{t,n}) \leq \varepsilon^2 \alpha. \end{aligned} \quad (5.22)$$

To prove (5.20) set  $C = \min\{\pi_0 V, \|\pi_0 - \pi\|_V\}$ , and combine (5.22) with (5.4) to get

$$n^2 - n \frac{\pi V \|f_c\|^p |V|^{2/p}}{\varepsilon^2 \alpha} \left(1 + \frac{2M_r \gamma_r}{1 - \gamma_r}\right) - \frac{MC \|f_c\|^p |V|^{2/p}}{\varepsilon^2 \alpha (1 - \gamma)} \left(1 + \frac{2M_r \gamma_r}{1 - \gamma_r}\right) \geq 0,$$

and hence  $n \geq \frac{b + \sqrt{b^2 + 4c}}{2}$ , where  $b$  and  $c$  are defined by (5.15) and (5.16) respectively. The only difference in (5.21) is that now we have  $c(t)$  defined

by (5.18) instead of  $c$ . It is easy to check that the best bound on  $t$  and  $n$  (i.e. that minimizes  $t + n$ ) is such that

$$n \geq n(t) \quad \text{and} \quad t \geq \max\{0, \min\{t \in N : n'(t) \geq -1\}\},$$

where  $n(t)$  is defined by (5.17). Standard calculations show that

$$\min\{t \in N : n'(t) \geq -1\} = \min\{t \in N : (\gamma^t)^2 \tilde{c}^2 \ln^2 \gamma - \gamma^t 4\tilde{c} - b^2 \leq 0\},$$

where  $\tilde{c}$  is defined by (5.19). Hence we obtain

$$t \geq \max\left\{0, (\ln \gamma)^{-1} \ln\left(\frac{2 + \sqrt{4 + b^2 \ln^2 \gamma}}{\tilde{c} \ln^2 \gamma}\right)\right\} \quad \text{and} \quad n \geq n(t).$$

This completes the proof.  $\square$

*Remark 5.4.2.* The formulation of Theorem 5.4.1 and the above proof indicate how the issue of a sufficient burn-in should be understood. The common description of  $t$  as *time to stationarity* and the often encountered approach that  $t^* = t(x, \tilde{\varepsilon})$  should be such that  $\rho(\pi, \delta_x P^{t^*}) \leq \tilde{\varepsilon}$  (where  $\rho(\cdot, \cdot)$  is a distance function for probability measures, e.g. total variation distance, or  $V$ -norm distance) seems not appropriate for such a natural goal as  $(\varepsilon - \alpha)$ -approximation. The optimal burn-in time can be much smaller than  $t^*$  and in particular cases it can be 0. Also we would like to emphasize that in the typical drift condition setting, i.e. if  $\mathcal{X}$  is not compact and the target function  $f$  is not bounded, the  $V$ -norm should be used as a measure of convergence, since  $\|\pi_t - \pi\|_{tv} \rightarrow 0$  does not even imply  $\pi_t f \rightarrow \pi f$ .

Next we suggest an alternative estimation scheme that allows for sharper bounds for the total simulation cost needed to obtain  $(\varepsilon - \alpha)$ -approximation for small  $\alpha$ . We will make use of the following simple lemma taken from the more complicated setting of [Niemiro & Pokarowski 2007].

**Lemma 5.4.3.** *Let  $m \in N$  be an odd number and let  $\hat{I}_1, \dots, \hat{I}_m$  be independent random variables, such that  $P(|\hat{I}_k - I| \leq \varepsilon) \geq 1 - a > 1/2$ , for  $k = 1, \dots, m$ . Define  $\hat{I} := \text{med}\{\hat{I}_1, \dots, \hat{I}_m\}$ . Then*

$$P(|\hat{I} - I| \leq \varepsilon) \geq 1 - \alpha, \quad \text{if} \quad m \geq \frac{2 \ln(2\alpha)}{\ln[4a(1-a)]}. \quad (5.23)$$



*Proof.* Since  $P(|\hat{I}_k - I| > \varepsilon) \leq a < 1/2$ , by elementary arguments we obtain

$$\begin{aligned} P(|\hat{I} - I| > \varepsilon) &\leq \sum_{k=(m+1)/2}^m \binom{m}{k} a^k (1-a)^{n-k} \\ &\leq 2^{m-1} a^{m/2} (1-a)^{m/2} \\ &= \frac{1}{2} \exp \left\{ \frac{m}{2} \ln(4a(1-a)) \right\}. \end{aligned}$$

The last term does not exceed  $\alpha$  if  $m \geq 2 \ln(2\alpha) / \ln[4a(1-a)]$ , as claimed.  $\square$

Hence  $(\varepsilon - \alpha)$ -approximation can be obtained by the following Algorithm 5.4.4, where Theorem 5.4.1 should be used to find  $t$  and  $n$  that guarantee  $(\varepsilon - a)$ -approximation and  $m$  results from Lemma 5.4.3.

**Algorithm 5.4.4.**

1. Simulate  $m$  independent runs of length  $t + n$  of the underlying Markov chain,

$$X_0^{(k)}, \dots, X_{t+n-1}^{(k)}, \quad k = 1, \dots, m.$$

2. Calculate  $m$  estimates of  $I$ , each based on a single run,

$$\hat{I}_k = \hat{I}_{t,n}^{(k)} = \frac{1}{n} \sum_{i=t}^{t+n-1} f(X_i^{(k)}), \quad k = 1, \dots, m.$$

3. For the final estimate take

$$\hat{I} = \text{med}\{\hat{I}_1, \dots, \hat{I}_m\}.$$

The total cost of Algorithm 5.4.4 amounts to

$$C = C(a) = m(t + n) \tag{5.24}$$

and depends on  $a$  (in addition to previous parameters). The optimal  $a$  can be found numerically, however it is worth mentioning  $a = 0,11969$  is an acceptable arbitrary choice (cf. [Niemiro & Pokarowski 2007]). A closer look at equation (5.24) reveals that the leading term is

$$mb = \frac{1}{a \ln\{[4a(1-a)]^{-1}\}} \left\{ \frac{2 \ln\{(2\alpha)^{-1}\} \pi V \|f_c\|^p |V|^{2/p}}{\varepsilon^2} \left( 1 + \frac{2M_r \gamma_r}{1 - \gamma_r} \right) \right\},$$

where  $b$  is defined by (5.15). Function  $a \ln\{[4a(1-a)]^{-1}\}$  has one maximum on  $(0, 1/2)$  at  $a \approx 0,11969$ .

## 5.5 A Toy Example - Contracting Normals

To illustrate the results of previous sections we analyze the *contracting normals* example studied by Baxendale in [Baxendale 2005] (see also [Roberts & Tweedie 1999], [Roberts & Rosenthal 1997a] and [Rosenthal 1995a]), where Markov chains with transition probabilities  $P(x, \cdot) = N(\theta x, 1 - \theta^2)$  for some parameter  $\theta \in (-1, 1)$  are considered.

Similarly as in [Baxendale 2005] we take a drift function  $V(x) = 1 + x^2$  and a small set  $C = [-c, c]$  with  $c > 1$ , which allows for  $\lambda = \theta^2 + \frac{2(1-\theta^2)}{1+c^2} < 1$  and  $K = 2 + \theta^2(c^2 - 1)$ . We also use the same minorization condition with  $\nu$  concentrated on  $C$ , such that  $\tilde{\beta}\nu(dy) = \min_{x \in C} (2\pi(1 - \theta^2))^{-1/2} \exp(-\frac{(\theta x - y)^2}{2(1-\theta^2)}) dy$ . This yields  $\tilde{\beta} = 2[\Phi(\frac{(1+|\theta|)c}{\sqrt{1-\theta^2}}) - \Phi(\frac{|\theta|c}{\sqrt{1-\theta^2}})]$ , where  $\Phi$  denotes the standard normal cumulative distribution function.

Baxendale in [Baxendale 2005] indicated that the chain is reversible with respect to its invariant distribution  $\pi = N(0, 1)$  for all  $\theta \in (-1, 1)$  and it is reversible and positive for  $\theta > 0$ .

Moreover, in Lemma 5.5.1 we observe a relationship between marginal distributions of the chain with positive and negative values of  $\theta$ . By  $\mathcal{L}(X_n|X_0, \theta)$  denote the distribution of  $X_n$  given the starting point  $X_0$  and the parameter value  $\theta$ .

**Lemma 5.5.1.**

$$\mathcal{L}(X_n|X_0, \theta) = \mathcal{L}(X_n|(-1)^n X_0, -\theta). \quad (5.25)$$

*Proof.* Let  $Z_1, Z_2, \dots$  be an iid  $N(0, 1)$  sequence, then

$$\begin{aligned} \mathcal{L}(X_n|X_0, \theta) &= \mathcal{L}\left(\theta^n X_0 + \sum_{k=1}^n \theta^{n-k} \sqrt{1 - \theta^2} Z_k\right) \\ &= \mathcal{L}\left((- \theta)^n (-1)^n X_0 + \sum_{k=1}^n (-\theta)^{n-k} \sqrt{1 - \theta^2} Z_k\right) \\ &= \mathcal{L}(X_n|(-1)^n X_0, -\theta), \end{aligned}$$

and we used the fact that  $Z_k$  and  $-Z_k$  have the same distribution.  $\square$

For  $\theta < 0$  using Lemma 5.5.1 and the fact that  $V(x) = 1 + x^2$  is symmetric we obtain

$$\begin{aligned} \|\mathcal{L}(X_n|X_0, \theta) - \pi\|_V &= \|\mathcal{L}(X_n|(-1)^n X_0, -\theta) - \pi\|_V \leq M\gamma^n V((-1)^n X_0) \\ &= M\gamma^n V(X_0) = M\gamma^n (1 + X_0^2). \end{aligned}$$

Thus for all  $\theta \in (-1, 1)$  we can bound the  $V$ -norm distance between  $\pi$  and the distribution of  $X_n$  via Theorem 5.2.4 with  $\rho$  and  $M = M(\gamma)$ , where  $\gamma \in (\rho, 1)$ , computed for reversible and positive Markov chains (see Appendix 5.8.3 for formulas). The choice of  $V(x) = 1 + x^2$  allows for

Table 5.5 - bounds based on Baxendale's  $V$ -uniform ergodicity constants.

$ \theta $	$\varepsilon$	$\alpha$	$\rho$	$\rho_2$	$\gamma$	$\gamma_2$	$M$	$M_2$	$m$	$t$	$n$	total cost
.5	.1	.1	.895	.899	.915	.971	36436	748	1	218	6.46e+09	6.46e+09
.5	.1	$10^{-5}$	.895	.899	.915	.971	36436	748	1	218	6.46e+13	6.46e+13
.5	.1	$10^{-5}$	.895	.899	.915	.971	36436	748	27	218	5.39e+09	1.46e+11

$(\varepsilon - \alpha)$ -approximation of  $\int_{\mathcal{X}} f(x)\pi(dx)$  if  $|f^2|_V < \infty$  for the possibly unbounded function  $f$ . In particular the MCMC works for all linear functions on  $\mathcal{X}$ . We take  $f(x) = x$  where  $|f^2|_V = 1$  as an example. We have to provide parameters and constants required for Theorem 5.4.1. In this case the optimal starting point is  $X_0 = 0$  since it minimizes  $V(x) = 1 + x^2$ . To bound  $\pi V$  we use Lemma 5.2.8 and Lemma 5.2.9 yields a bound on  $\|f_c\|_V^{2/p} = |f_c^2|_V$ .

Examples of bounds for  $t$  and  $n$  for the one walk estimator, or  $t$ ,  $n$  and  $m$  for the median of multiple runs estimator are given in Table 5.5. The bounds are computed for  $c = 1.6226$  which minimizes  $\rho_2$  (rather than  $\rho$ ) for  $\theta = 0.5$ . Then a grid search is performed to find optimal values of  $\gamma$  and  $\gamma_2$  that minimize the total simulation cost. Note that in Baxendale's constant  $M$  depends on  $\gamma$  and  $M$  goes relatively quickly to  $\infty$  as  $\gamma \rightarrow \rho$ . This is the reason why optimal  $\gamma$  and  $\gamma_2$  are far from  $\rho$  and  $\rho_2$  and this turns out to be the main weakness of Baxendale's bounds. Also for small  $\alpha = 10^{-5}$  we observe a clear computational advantage of the median of multiple runs estimation. The  $m = 27$  shorter runs have significantly lower total cost than the single long run. R functions for computing this example and also the general bounds resulting from Theorem 5.4.1 are available at <http://akson.sgh.waw.pl/~klatus/>

## 5.6 The Example - a Hierarchical Random Effects Model

In this section we describe a hierarchical random effects model which is a widely applicable example and provides a typical target density  $\pi$  that arises

in Bayesian statistics. Versions of this model and the efficiency of MCMC sampling have been analyzed e.g. by Gelfand and Smith in [Gelfand & Smith 1990], Rosenthal in [Rosenthal 1995a], [Rosenthal 1995b] and many other authors. In particular Hobert and Geyer in [Hobert & Geyer 1998] analyzed a Gibbs sampler and a block Gibbs sampler for this model and showed the underlying Markov chains are in both cases geometrically ergodic (we describe these samplers in the sequel). Jones and Hobert in [Jones & Hobert 2004] derived computable bounds for the geometric ergodicity parameters and consequently computable bounds for the total variation distance  $\|P^t(x, \cdot) - \pi\|_{tv}$  to stationarity in both cases. They used these bounds to determine the burn-in time. Their work was a breakthrough in analyzing the hierarchical random effects model, however, mere bounds on burn-in time do not give a clue on the total amount of simulation needed. Also, bounding the total variation distance seems inappropriate when estimating integrals of unbounded functions, as indicated in Remark 5.4.2. In this section we establish the  $(\varepsilon - \alpha)$ -approximation for the hierarchical random effects model. This consists of choosing a suitable sampler, establishing the Drift Condition 5.2.1 with explicit constants, computing  $V$ -uniform ergodicity parameters, and optimizing lower bounds for  $t$  and  $n$  in case of estimation along one walk or for  $t$ ,  $n$  and  $m$  in (5.24) for the median of shorter runs. This may turn out to be a confusing procedure, hence we outline it here in detail, discuss computational issues and provide necessary R functions.

### 5.6.1 The Model

Since we will make use of the drift conditions established by Jones and Hobert in [Jones & Hobert 2004] we also try to follow their notation in the model description. Let  $\mu$  and  $\lambda_\theta$  be independent and distributed as

$$\mu \sim N(m_0, s_0^{-1}) \quad \text{and} \quad \lambda_\theta \sim \text{Gamma}(a_1, b_1),$$

where  $m_0 \in \mathbb{R}$ ,  $s_0 > 0$ ,  $a_1 > 0$ , and  $b_1 > 0$  are known constants.

At the second stage, conditional on  $\mu$  and  $\lambda_\theta$ , random variables  $\theta_1, \dots, \theta_K$  and  $\lambda_e$  are independent and distributed as

$$\theta_i | \mu, \lambda_\theta \sim N(\mu, \lambda_\theta^{-1}) \quad \text{and} \quad \lambda_e \sim \text{Gamma}(a_2, b_2),$$

where  $a_2 > 0$ ,  $b_2 > 0$  are known constants.

Finally in the third stage, conditional on  $\theta = (\theta_1, \dots, \theta_K)$  and  $\lambda_e$ , the observed data  $y = \{Y_{ij}\}$  are independent with

$$Y_{ij}|\theta, \lambda_e \sim N(\theta_i, \lambda_e^{-1}),$$

where  $i = 1, \dots, K$  and  $j = 1, \dots, m_i$ .

The Bayesian approach involves conditioning on the values of the observed data  $\{Y_{ij}\}$  and considering the joint distribution of all  $K+3$  parameters given this data. Thus we are interested in the posterior distribution, that is, the following distribution defined on the space  $\mathcal{X} = (0, \infty)^2 \times \mathbb{R}^{K+1}$ ,

$$\begin{aligned} \mathcal{L}(\theta_1, \dots, \theta_K, \mu, \lambda_\theta, \lambda_e | \{Y_{ij}\}) &= \pi(\theta, \mu, \lambda | y) \\ &\propto d(y|\theta, \lambda_e) d(\theta|\mu, \lambda_\theta) d(\lambda_e) d(\lambda_\theta) d(\mu) = \clubsuit, \end{aligned} \quad (5.26)$$

where  $d$  denotes a generic density and hence the final formula for the unnormalised density takes the form of

$$\begin{aligned} \clubsuit &= e^{-b_1 \lambda_\theta} \lambda_\theta^{a_1-1} e^{-b_2 \lambda_e} \lambda_e^{a_2-1} e^{-\frac{1}{2} s_0 (\mu - m_0)^2} \\ &\times \prod_{i=1}^K \left[ e^{-\frac{1}{2} \lambda_\theta (\theta_i - \mu)^2} \lambda_\theta^{1/2} \right] \times \prod_{i=1}^K \prod_{j=1}^{m_i} \left[ e^{-\frac{1}{2} \lambda_e (y_{ij} - \theta_i)^2} \lambda_e^{1/2} \right], \end{aligned} \quad (5.27)$$

and we have to deal with a density that is high-dimensional, irregular, strictly positive in  $\mathcal{X}$  and concentrated in the „center” of  $\mathcal{X}$ , which is very typical for MCMC situations [Roberts & Rosenthal 2005]. Computing expectations with respect to  $\pi(\theta, \mu, \lambda | y)$  is crucial for bayesian inference (e.g. to obtain bayesian estimators) and requires MCMC techniques.

## 5.6.2 Gibbs Samplers for the Model

Full conditional distributions required for a Gibbs sampler can be computed without difficulty. Let

$$\begin{aligned} \bar{y}_i &:= \frac{1}{m_i} \sum_{j=1}^{m_i} y_{ij}, & M &:= \sum_i m_i, & \bar{\theta} &= \frac{1}{K} \sum_i \theta_i, \\ \theta_{-i} &:= (\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_K), & \nu_1(\theta, \mu) &:= \sum_{i=1}^K (\theta_i - \mu)^2, \end{aligned}$$

$$\nu_2(\theta) := \sum_{i=1}^K (\theta_i - \bar{y}_i)^2, \quad SSE := (y_{ij} - \bar{y}_i)^2.$$

Now the conditionals are

$$\lambda_\theta | \theta, \mu, \lambda_e, y \sim \text{Gamma} \left( \frac{K}{2} + a_1, \frac{\nu_1(\theta, \mu)}{2} + b_1 \right), \quad (5.28)$$

$$\lambda_e | \theta, \mu, \lambda_\theta, y \sim \text{Gamma} \left( \frac{M}{2} + a_2, \frac{\nu_2(\theta) + SSE}{2} + b_2 \right), \quad (5.29)$$

$$\theta_i | \theta_{-i}, \mu, \lambda_\theta, \lambda_e, y \sim N \left( \frac{\lambda_\theta \mu + m_i \lambda_e \bar{y}_i}{\lambda_\theta + m_i \lambda_e}, \frac{1}{\lambda_\theta + m_i \lambda_e} \right), \quad (5.30)$$

$$\mu | \theta, \lambda_\theta, \lambda_e, y \sim N \left( \frac{s_0 m_0 + K \lambda_\theta \bar{\theta}}{s_0 + K \lambda_\theta}, \frac{1}{s_0 + K \lambda_\theta} \right). \quad (5.31)$$

Gibbs samplers for the variance components model and its versions have been used and studied by many authors. We consider the two Gibbs samplers analyzed by Jones and Hobert in [Jones & Hobert 2004].

- The **fixed-scan Gibbs sampler** that updates  $\mu$ , then  $\theta = (\theta_1, \dots, \theta_K)$ , then  $\lambda = (\lambda_\theta, \lambda_e)$ . Note that  $\theta_i$ 's are conditionally independent given  $(\mu, \lambda)$  and so are  $\lambda_\theta$  and  $\lambda_e$  given  $(\theta, \mu)$ . Thus the one step Markov transition density  $(\mu', \theta', \lambda') \rightarrow (\mu, \theta, \lambda)$  of this Gibbs sampler is

$$p(\mu, \theta, \lambda | \mu', \theta', \lambda') = d(\mu | \theta', \lambda', y) \left[ \prod_{i=1}^K d(\theta_i | \mu, \lambda', y) \right] \times d(\lambda_\theta | \theta, \mu, y) d(\lambda_e | \theta, \mu, y). \quad (5.32)$$

Where  $d$  denotes a generic density and  $y = \{Y_{ij}\}$ ,  $i = 1, \dots, K$ ;  $j = 1, \dots, m_i$ , is the observed data.

- Hobert and Geyer in [Hobert & Geyer 1998] introduced a more efficient **block Gibbs sampler** (also analyzed by Jones and Hobert in [Jones & Hobert 2004]), in which all the components of

$$\xi = (\theta_1, \dots, \theta_K, \mu)$$

are updated simultaneously. It turns out that

$$\xi | \lambda, y \sim N(\xi^*, \Sigma) \quad \text{where} \quad \xi^* = \xi^*(\lambda, y) \quad \text{and} \quad \Sigma = \Sigma(\lambda, y).$$

Thus the one step Markov transition density  $(\lambda', \xi') \rightarrow (\lambda, \xi)$  of the block Gibbs sampler is

$$p(\lambda, \xi | \lambda', \xi') = d(\lambda_\theta | \xi', y) d(\lambda_e | \xi', y) d(\xi | \lambda, y). \quad (5.33)$$

We give now the formulas for  $\xi^*$  and  $\Sigma$  derived in [Hobert & Geyer 1998]. Let

$$\tau = \sum_{i=1}^K \frac{m_i \lambda_\theta \lambda_e}{\lambda_\theta + m_i \lambda_e},$$

then

$$\begin{aligned} E(\mu | \lambda) &= \frac{1}{s_0 + \tau} \left[ \sum_{i=1}^K \frac{m_i \lambda_\theta \lambda_e \bar{y}_i}{\lambda_\theta + m_i \lambda_e} + m_0 s_0 \right], \\ E(\theta_i | \lambda) &= \frac{\lambda_\theta E(\mu | \lambda)}{\lambda_\theta + m_i \lambda_e} + \frac{m_i \lambda_e \bar{y}_i}{\lambda_\theta + m_i \lambda_e}. \end{aligned}$$

and

$$\begin{aligned} Var(\theta_i | \lambda) &= \frac{1}{\lambda_\theta + m_i \lambda_e} \left[ 1 + \frac{\lambda_\theta^2}{(\lambda_\theta + m_i \lambda_e)(s_0 + \tau)} \right], \\ Cov(\theta_i, \theta_j | \lambda) &= \frac{\lambda_\theta^2}{\lambda_\theta + m_i \lambda_e (\lambda_\theta + m_j \lambda_e) (s_0 + \tau)}, \\ Cov(\theta_i, \theta_j | \lambda) &= \frac{\lambda_\theta^2}{\lambda_\theta + m_i \lambda_e (\lambda_\theta + m_j \lambda_e) (s_0 + \tau)}, \\ Var(\mu | \lambda) &= \frac{1}{s_0 + \tau}. \end{aligned}$$

### 5.6.3 Relations between Drift Conditions

A crucial step for  $(\varepsilon - \alpha)$ -approximation is establishing the drift condition 5.2.1 which in the sequel will be referred to as the Baxendale-type drift condition. To this end we use the Rosenthal-type (cf. [Rosenthal 1995b]) and Roberts-and-Tweedie-type (cf. [Roberts & Tweedie 1999]) drift conditions established by Jones and Hobert in [Jones & Hobert 2004] combined with their type of a small set condition.

In the following definitions and lemmas  $P$  denotes the transition kernel of the Markov chain  $(X_n)_{n \geq 0}$  and the subscripts of drift condition parameters indicate the type of drift condition they refer to.

**Assumption 5.6.1** (The Rosenthal-type drift condition).

(R.1) *There exists a function  $V_R : \mathcal{X} \rightarrow [0, \infty)$  and constants  $0 < \lambda_R < 1$  and  $K_R < \infty$  satisfying*

$$PV_R(x) \leq \lambda_R V_R(x) + K_R. \quad (5.34)$$

(R.2) *Let  $C_R = \{x \in \mathcal{X} : V_R(x) \leq d_R\}$ , where  $d_R > 2K_R/(1 - \lambda_R)$ . There exists a probability measure  $\nu_R$  on  $\mathcal{X}$  and  $\tilde{\beta}_R > 0$ , such that for all  $x \in C_R$  and  $A \in \mathcal{B}(\mathcal{X})$ ,*

$$P(x, A) \geq \tilde{\beta}_R \nu_R(A). \quad (5.35)$$

**Assumption 5.6.2** (The Roberts-and-Tweedie-type drift condition).

(RT.1) *There exists a function  $V_{RT} : \mathcal{X} \rightarrow [1, \infty)$  and constants  $0 < \lambda_{RT} < 1$  and  $K_{RT} < \infty$  satisfying*

$$PV_{RT}(x) \leq \lambda_{RT} V_{RT}(x) + K_{RT} \mathbb{I}_{C_{RT}}(x), \quad (5.36)$$

where  $C_{RT} = \{x \in \mathcal{X} : V_{RT}(x) \leq d_{RT}\}$ , and  $d_{RT} \geq \frac{K_{RT}}{1 - \lambda_{RT}} - 1$ .

(RT.2) *There exists a probability measure  $\nu_{RT}$  on  $\mathcal{X}$  and  $\tilde{\beta}_{RT} > 0$ , such that for all  $x \in C_{RT}$  and  $A \in \mathcal{B}(\mathcal{X})$ ,*

$$P(x, A) \geq \tilde{\beta}_{RT} \nu_{RT}(A). \quad (5.37)$$

The following lemma relates the two drift conditions.

**Lemma 5.6.3** (Lemma 3.1 of [Jones & Hobert 2004]). *Assume that the Rosenthal-type drift condition holds. Then for any  $d > 0$  the Roberts-and-Tweedie-type drift condition holds with parameters*

$$V_{RT} = V_R + 1, \quad \lambda_{RT} = \lambda_{RT}(d) = \frac{d + \lambda_R}{d + 1}, \quad K_{RT} = K_R + 1 - \lambda_R, \quad \tilde{\beta}_{RT} = \tilde{\beta}_R,$$

$$C_{RT} = C_{RT}(d) = \left\{ x \in \mathcal{X} : V_{RT}(x) \leq \frac{(d + 1)K_{RT}}{d(1 - \lambda_{RT})} \right\} \quad \text{and} \quad \nu_{RT} = \nu_R.$$

The Baxendale-type drift condition we work with results from each of the above conditions and the following lemma is easy to verify by simple algebra.



**Lemma 5.6.4.** *If the Rosenthal-type or the Roberts-and-Tweedie-type drift condition holds, then the Baxendale-type drift condition (A.1-2) verifies with*

$$\begin{aligned}
V &= V_{RT} = V_R + 1, & \lambda &= \lambda(d) = \lambda_{RT} = \frac{d + \lambda_R}{d + 1}, \\
\nu &= \nu_{RT} = \nu_R, & C &= C(d) = C_{RT}, & \tilde{\beta} &= \tilde{\beta}_{RT} = \tilde{\beta}_R, \\
K &= K(d) = K_{RT} + \lambda_{RT}d_{RT} = (K_R + 1 - \lambda_R) \frac{d^2 + 2d + \lambda_R}{d(1 - \lambda_R)}.
\end{aligned}$$

Observe next that integrating each of the drift conditions yields a bound on  $\pi V$  similar to the one obtained in Lemma 5.2.8 and the best available bound should be used in Theorem 5.3.4 and Theorem 5.4.1. In particular, if the Baxendale-type drift condition is obtained from the Roberts-and-Tweedie-type drift condition via Lemma 5.6.4, integrating the latter always leads to a better bound on  $\pi V$ . Also, if one starts with establishing the Rosenthal-type drift condition, the value of  $d$  used for bounding  $\pi V$  does not have to be the same as the one used for establishing the Baxendale-type drift and minorization condition and it should be optimized. Moreover  $\frac{K_R}{1 - \lambda_R} + 1 < \frac{K_{RT}}{1 - \lambda_{RT}} < \frac{K - \lambda}{1 - \lambda}$  for every  $d > 0$ . This leads to the following lemma which can be checked by straightforward calculations.

**Lemma 5.6.5.** *Provided the drift functions are as in Lemma 5.6.4, the bound on  $\pi V$  can be optimized as follows*

$$\pi V \leq \min \left\{ \inf_d \left\{ \pi(C_{RT}(d)) \frac{K_{RT}}{1 - \lambda_{RT}(d)} \right\}, \frac{K_R}{1 - \lambda_R} + 1 \right\} \leq \frac{K_R}{1 - \lambda_R} + 1. \tag{5.38}$$

### 5.6.4 Drift and Minorization Conditions for the Samplers

For the fixed-scan Gibbs sampler and the block Gibbs sampler of Section 5.6.2 Jones and Hobert in [Jones & Hobert 2004] (Section 4 and 5 therein) obtained the following drift and minorization conditions. See their paper for derivation and more elaborative commentary of these results.

### Drift and Minorization for the block Gibbs Sampler

Assume  $m' = \min\{m_1, \dots, m_K\} \geq 2$  and  $K \geq 3$ . Moreover define

$$\delta_1 = \frac{1}{2a_1 + K - 2}, \quad \delta_2 = \frac{1}{2a_2 + M - 2}, \quad \delta_3 = (K + 1)\delta_2, \quad \delta_4 = \delta_2 \sum_{i=1}^K m_i^{-1},$$

$$\delta = \max\{\delta_1, \delta_3\}, \quad c_1 = \frac{2b_1}{2a_1 + K - 2}, \quad c_2 = \frac{2b_2 + SSE}{2a_2 + M - 2}.$$

Observe that  $0 < \delta_i < 1$  for  $i = 1, 2, 3, 4$ . Also let  $\Delta$  denote the length of the convex hull of the set  $\{\bar{y}_1, \dots, \bar{y}_K, m_0\}$ .

**Proposition 5.6.6** (Drift for unbalanced case). *Fix  $\lambda_R \in (\delta, 1)$  and let  $\phi_1$  and  $\phi_2$  be positive numbers such that  $\frac{\phi_1 \delta_4}{\phi_2} + \delta < \lambda_R$ . Define the drift function as*

$$V_1(\theta, \mu) = \phi_1 \nu_1(\theta, \mu) + \phi_2 \nu_2(\theta), \quad (5.39)$$

where  $\nu_1(\theta, \mu)$  and  $\nu_2(\theta)$  are defined in Section 5.6.2. With this drift function the block Gibbs sampler satisfies the Rosenthal-type drift condition with

$$K_R = \phi_1 \left[ c_1 + c_2 \frac{\delta_4}{\delta_2} + K \Delta^2 \right] + \phi_2 \left[ c_2 (K + 1) + M \Delta^2 \right]. \quad (5.40)$$

A better drift condition can be obtained in the balanced case, when  $m_i = m \geq 2$  for  $i = 1, \dots, K$ . Let  $\delta_5 = K \delta_2$ .

**Proposition 5.6.7** (Drift for balanced case). *Fix  $\lambda_R \in (\delta, 1)$  and let  $\phi$  be a positive number such that  $\phi \delta_5 + \delta < \lambda_R$ . Define the drift function as*

$$V_2(\theta, \mu) = \phi \nu_1(\theta, \mu) + m^{-1} \nu_2(\theta). \quad (5.41)$$

With this drift function the block Gibbs sampler satisfies the Rosenthal-type drift condition with

$$K_R = \phi c_1 + (\phi K + K + 1) \frac{c_2}{m} + \max\{\phi, 1\} \sum_{i=1}^K \max\{(\bar{y} - \bar{y}_i)^2, (m_0 - \bar{y}_i)^2\}, \quad (5.42)$$

where  $\bar{y} := K^{-1} \sum_{i=1}^K \bar{y}_i$ .

Proposition 5.6.8 (Proposition 4.1 of [Jones & Hobert 2004]) provides a minorization condition for the Rosenthal-type drift-minorization condition for the block Gibbs sampler for both, the balanced and unbalanced case. Note that the balanced case drift function  $V_2$  is a special case of the unbalanced drift function  $V_1$ , hence we focus on  $V_1$ .

Now consider the candidate  $C_R = \{(\theta, \mu) : V_1(\theta, \mu) \leq d_R\}$  for a small set. Note that  $C_R$  is contained in  $S_B = S_{B_1} \cap S_{B_2}$ , where  $S_{B_1} = \{(\theta, \mu) : \nu_1(\theta, \mu) < d_R/\phi_1\}$  and  $S_{B_2} = \{(\theta, \mu) : \nu_2(\theta) < d_R/\phi_2\}$ . Hence it is enough to establish a minorization condition that holds for  $S_B$ .

Let  $\Gamma(\alpha, \beta; x)$  denote the value of the Gamma( $\alpha, \beta$ ) density at  $x$  and define functions  $h_1(\lambda_\theta)$  and  $h_2(\lambda_e)$  as follows:

$$h_1(\lambda_\theta) = \begin{cases} \Gamma\left(\frac{K}{2} + a_1, b_1; \lambda_\theta\right), & \lambda_\theta < \lambda_\theta^*, \\ \Gamma\left(\frac{K}{2} + a_1, \frac{d_R}{2\phi_1} + b_1; \lambda_\theta\right), & \lambda_\theta \geq \lambda_\theta^*, \end{cases}$$

where

$$\lambda_\theta^* = \frac{\phi_1(K + 2a_1)}{d_R} \log\left(1 + \frac{d_R}{2b_1\phi_1}\right)$$

and

$$h_2(\lambda_e) = \begin{cases} \Gamma\left(\frac{M}{2} + a_2, \frac{SSE}{2} + b_2; \lambda_e\right), & \lambda_e < \lambda_e^*, \\ \Gamma\left(\frac{M}{2} + a_2, \frac{d_R + \phi_2 SSE}{2\phi_2} + b_2; \lambda_e\right), & \lambda_e \geq \lambda_e^*, \end{cases}$$

where

$$\lambda_e^* = \frac{\phi_2(M + 2a_2)}{d_R} \log\left(1 + \frac{d_R}{\phi_2(2b_2 + SSE)}\right).$$

Now define a density  $q(\lambda, \theta, \mu)$  on  $R_+^2 \times R^K \times R$  by

$$q(\lambda, \theta, \mu) = \left(\frac{h_1(\lambda_\theta)}{\int_{R_+} h_1(\lambda_\theta) d\lambda_\theta}\right) \left(\frac{h_2(\lambda_e)}{\int_{R_+} h_2(\lambda_e) d\lambda_e}\right) d(\xi|\lambda, y),$$

where  $d(\xi|\lambda, y)$  is the normal density in (5.33) resulting from the block Gibbs sampler construction. Next define

$$\tilde{\beta}_R = \left(\int_{R_+} h_1(\lambda_\theta) d\lambda_\theta\right) \left(\int_{R_+} h_2(\lambda_e) d\lambda_e\right).$$

Also recall  $p(\lambda, \xi|\lambda', \xi') = p(\lambda, \theta, \mu|\lambda', \theta', \mu')$ , the Markov transition density of the block Gibbs sampler as specified in (5.33).

We are in a position to state the minorization condition.

**Proposition 5.6.8** (Minorization Condition). *The Markov transition density for the block Gibbs sampler satisfies the following minorization condition:*

$$p(\lambda, \theta, \mu | \lambda', \theta', \mu') \geq \tilde{\beta}_R q(\lambda, \theta, \mu) \quad \text{for every } (\theta', \mu') \in S_B. \quad (5.43)$$

### Drift and Minorization for the fixed-scan Gibbs sampler

As before assume that  $K \geq 3$  and

$$2 \leq m' = \min\{m_1, \dots, m_K\} \leq \max\{m_1, \dots, m_K\} = m''.$$

Define

$$\delta_6 = \frac{K^2 + 2Ka_1}{2s_0 + K^2 + 2Ka_1} \quad \text{and} \quad \delta_7 = \frac{1}{2(a_1 - 1)}.$$

Clearly  $\delta_6 \in (0, 1)$  and if  $a_1 > 3/2$  then also  $\delta_7 \in (0, 1)$ . Moreover if  $a_1 > 3/2$ , then since  $2s_0b_1 > 0$ , there exists  $\rho_1 \in (0, 1)$  such that

$$\left(K + \frac{\delta_6}{\delta_7}\right)\delta_1 < \rho_1. \quad (5.44)$$

Define also

$$\nu_3(\theta, \lambda) = \frac{K\lambda_\theta}{s_0 + K\lambda_\theta}(\bar{\theta} - \bar{y})^2 \quad \text{and} \quad s^2 = \sum_{i=1}^K (\bar{y}_i - \bar{y})^2.$$

**Proposition 5.6.9** (Drift Condition). *Assume that  $a_1 > 3/2$ ,  $5m' > m''$  and let  $\rho_1 \in (0, 1)$  satisfy (5.44). Fix*

$$c_3 \in (0, \min\{b_1, b_2\}) \quad \text{and} \quad \lambda_R \in (\max\{\rho_1, \delta_6, \delta_7\}, 1).$$

Define the drift function as

$$V_3(\theta, \lambda) = e^{c_3\lambda_\theta} + e^{c_3\lambda_e} + \frac{\delta_7}{K\delta_1\lambda_\theta} + \nu_3(\theta, \lambda). \quad (5.45)$$

With this drift function the fixed-scan Gibbs sampler satisfies the Rosenthal-type drift condition with

$$K_R = \left(\frac{b_1}{b_1 - c_3}\right)^{a_1 + \frac{K}{2}} + \left(\frac{b_2}{b_2 - c_3}\right)^{a_2 + \frac{M^2}{2}} + (\delta_6 + \delta_7) \left[\frac{1}{s_0} + (m_0 - \bar{y})^2 + \frac{s^2}{K}\right] + \frac{2b_1\delta_7}{K}. \quad (5.46)$$

We now turn to the minorization condition for the fixed-scan Gibbs sampler provided in Section 5.2 of [Jones & Hobert 2004]. Similarly as before, consider the candidate  $C_R = \{(\theta, \lambda) : V_3 \leq d_R\}$  for a small set and let

$$c_4 = \frac{\delta_7}{K\delta_1 d_R}, \quad c_l = \bar{y} - \sqrt{(m_0 - \bar{y})^2 + d_R} \quad \text{and} \quad c_u = \bar{y} + \sqrt{(m_0 - \bar{y})^2 + d_R}.$$

The minorization condition will be given on a set  $S_G$  such that

$$C_R \subseteq S_G = S_{G_1} \cap S_{G_2} \cap S_{G_3},$$

where

$$\begin{aligned} S_{G_1} &= \left\{ (\theta, \lambda) : c_4 \leq \lambda_\theta \leq \frac{\log d_R}{c_3} \right\}, \\ S_{G_2} &= \left\{ (\theta, \lambda) : 0 < \lambda_e \leq \frac{\log d_R}{c_3} \right\}, \\ S_{G_3} &= \left\{ (\theta, \lambda) : c_l \leq \frac{s_0 m_0 + K \lambda_\theta \bar{\theta}}{s_0 + K \lambda_\theta} \leq c_u \right\}. \end{aligned}$$

Moreover to assure that  $S_{G_1} \cap S_{G_2}$  is nonempty, choose  $d_R$  such that

$$d_R \log d_R > \frac{c_3 \delta_7}{K \delta_1}.$$

Let  $N(\zeta, \sigma^2; x)$  denote the value of the  $N(\zeta, \sigma^2)$  density at  $x$  and define functions  $g_1(\mu, \theta)$  and  $g_2(\mu)$  as follows:

$$g_1(\nu, \theta) = \left( \frac{c_4}{2\pi} \right)^{\frac{K}{2}} \exp \left\{ - \frac{\log d_R}{2c_3} \sum_{i=1}^K [(\theta_i - \mu)^2 + m_i(\theta_i - \bar{y}_i)^2] \right\}$$

and

$$g_2(\mu) = \begin{cases} N(c_u, [s_0 + \frac{K \log(d_R)}{c_3}]^{-1}; \mu), & \mu \leq \bar{y}, \\ N(c_l, [s_0 + \frac{K \log(d_R)}{c_3}]^{-1}; \mu), & \mu > \bar{y}. \end{cases}$$

Now define a density on  $R \times R^K \times R_+^2$  by

$$q(\mu, \theta, \lambda) = \left( \frac{g_1(\mu, \theta) g_2(\mu)}{\int_R \int_{R^K} g_1(\mu, \theta) g_2(\mu) d\theta d\mu} \right) d(\lambda | \mu, \theta, y),$$

where  $d(\lambda|\mu, \theta, y)$  is the joint Gamma distribution of  $\lambda_\theta$  and  $\lambda_e$  in (5.32) resulting from the fixed-scan Gibbs sampler construction. Next define

$$\tilde{\beta}_R = \left( \frac{s_0 + Kc_4}{s_0 + \frac{K \log d_R}{c_3}} \right)^{1/2} \left( \int_R \int_{R^K} g_1(\mu, \theta) g_2(\mu) d\theta d\mu \right).$$

Also recall  $p(\mu, \theta, \lambda|\mu', \theta', \lambda')$ , the Markov transition density of the fixed-scan Gibbs sampler as specified in (5.32). We are in a position to state the minorization condition.

**Proposition 5.6.10.** *The Markov transition density for the fixed-scan Gibbs sampler satisfies the following minorization condition*

$$p(\mu, \theta, \lambda|\mu', \theta', \lambda') \geq \tilde{\beta}_R q(\mu, \theta, \lambda) \quad \text{for every } (\theta', \lambda') \in S_G. \quad (5.47)$$

Moreover Jones and Hobert in [Jones & Hobert 2004] obtained a closed form expression for  $\tilde{\beta}_R$  in (5.47) involving the standard normal cumulative distribution function  $\Phi$ . Let

$$\begin{aligned} \nu &= \left[ s_0 + \frac{\log d_R}{c_3} \left( K + \sum_{i=1}^K \frac{m_i}{1 + m_i} \right) \right]^{-1}, \\ m_l &= \nu \left[ c_l s_0 + \frac{\log d_R}{c_3} \left( K c_l + \sum_{i=1}^K \frac{\bar{y}_i m_i}{1 + m_i} \right) \right], \\ m_u &= \nu \left[ c_u s_0 + \frac{\log d_R}{c_3} \left( K c_u + \sum_{i=1}^K \frac{\bar{y}_i m_i}{1 + m_i} \right) \right]. \end{aligned}$$

Then

$$\begin{aligned} \tilde{\beta}_R &= \left( \frac{c_4 c_3}{\log d_R} \right)^{\frac{K}{2}} \sqrt{\nu (s_0 + K c_4)} \sqrt{\prod_{i=1}^K \frac{1}{1 + m_i} \exp \left\{ -\frac{\log d_R}{2c_3} \sum_{i=1}^K \frac{\bar{y}_i^2 m_i}{1 + m_i} \right\}} \\ &\times \left[ \exp \left\{ -\frac{c_u^2 s_0}{2} - \frac{K c_u^2 \log d_R}{2c_3} + \frac{m_u^2}{2\nu} \right\} \Phi \left( \frac{\bar{y} - m_u}{\sqrt{\nu}} \right) \right. \\ &\quad \left. + \exp \left\{ -\frac{c_l^2 s_0}{2} - \frac{K c_l^2 \log d_R}{2c_3} + \frac{m_l^2}{2\nu} \right\} \left( 1 - \Phi \left( \frac{\bar{y} - m_l}{\sqrt{\nu}} \right) \right) \right]. \end{aligned}$$

### 5.6.5 Obtaining the Bounds

We focus on obtaining the bounds for  $(\varepsilon - \alpha)$ -approximation for bayesian estimators of parameters  $\mu, \lambda_\theta, \lambda_e$  and  $\theta_i$ . This involves integrating one dimensional projections of the identity function on parameter space. The drift function  $V$  has to be at least of order  $f^2$  since  $|f^2|_V$  has to be finite. Note that for the two described samplers different drift conditions has been established and neither of them majorizes quadratic functions in all the parameters. Thus specifying a parameter, say  $\lambda_e$  implies the choice of the fixed-scan Gibbs sampler with the drift function  $V_3$ , whereas for  $\mu$  the block-scan Gibbs sampler with drift function  $V_1$  or  $V_2$  is the only option.

Once the sampler and the type of the drift condition is chosen, the user must provide his choice of  $\lambda_R, \phi$  and  $d_R$  for the Rosenthal-type drift-minorization condition. The next step is the right choice of  $d$  in Lemma 5.6.4 which yields the parameters of the Baxendale-type drift condition. Provided the Baxendale-type drift condition is established with computable parameters, there are still four parameters left to the user, namely the mutually dependent  $\gamma$  and  $M$  in Baxendale's Theorem 5.2.4 and their counterparts  $\gamma_2$  and  $M_2$  from Corollary 5.2.7. Unfortunately the bounds on  $t$  and  $n$  or  $t, n$  and  $m$  are very complicated functions of these parameters subject to users choice and finding optimal values analytically seems impossible. Also, in our experience, small changes in these quantities usually result in dramatically different bounds.

Similarly as burn-in bounds in [Jones & Hobert 2004], final bounds for  $(\varepsilon - \alpha)$ -approximation also strongly depend on the hyperparameter setting and the observed data.

Thus we provide appropriate R functions for approximating optimal bonds on the simulation parameters. This functions are available on <http://akson.sgh.waw.pl/~klatus/>

## 5.7 Concluding Remarks

To our best knowledge, in the above setting of an unbounded target function  $f$  and without assuming uniform ergodicity of the underlying Markov chain (which in practice means the state space  $\mathcal{X}$  is not compact) we derived first explicit bounds for the total simulation cost required for  $(\varepsilon - \alpha)$ -approximation. These bounds are sometimes feasible and sometimes infeasible on a PC, and probably always exceed the true values by many orders

of magnitude. Although  $10^9$  iterations in our Toy Example takes about 1 minute on a standard PC, sampling more realistic chains will take more time and the bound will be even more conservative.

However, the message of the Chapter is a very positive one: the current theoretical knowledge of Markov chains has reached the stage when for many MCMC algorithms of practical relevance applied to difficult problems, i.e. estimating expectations of unbounded functions, we are able to provide a rigorous, nonasymptotic, a priori analysis of the quality of estimation. This is much more than the often used in practice visual assessment of convergence, more sophisticated a posteriori convergence diagnostics, bounding only burn in time or even using asymptotic confidence intervals.

We also notice the following:

- The leading term in the bound for  $n$  is  $b = \frac{\pi V|f_c^2|_V}{\varepsilon^2 \alpha} (1 + \frac{2M_2\gamma_2}{1-\gamma_2})$  (where we took  $p = r = 2$  for simplicity).  $\pi V|f_c^2|_V$  should be of the order of  $Var_{\pi} f$ , thus this term is inevitable.  $\varepsilon^{-2}$  results from Chebyshev's inequality, since we proceed by bounding the mean square error.  $\alpha^{-1}$  can be reduced to  $\log(\alpha^{-1})$  for small  $\alpha$  by Lemma 5.4.3 and Algorithm 5.4.4 which in fact results in an exponential inequality. The last term  $1 + \frac{2M_2\gamma_2}{1-\gamma_2}$  is of the same order as a general bound for the ratio of the asymptotic variance and the stationary variance, under drift condition and without reversibility as indicated by Theorem 5.3.4. Thus it also seems to be inevitable. However we acknowledge this bound seems to be very poor due to the present form of  $V$ -uniform ergodicity constants.
- The term  $1 + \frac{2M_2\gamma_2}{1-\gamma_2}$  is the bottleneck of the approach. Here good bounds on  $\gamma$  and the somewhat disregarded  $M(\gamma)$  are equally important. Improvements in Baxendale-type convergence bounds may lead to dramatic improvement of the bounds on the total simulation cost (e.g. by applying the preliminary results of [Bednorz 2008]).
- Improvements of drift parameters (i.e. establishing better drift functions and minorization conditions) imply significant improvement of the convergence bounds in Baxendale's Theorem.
- The drift conditions we used as well as the Baxendale's theorem are far from optimal and subject to improvement.
- We applied the theoretical results to the toy example of Section 5.5 where the drift and minorization conditions are available without much



effort and to the Hierarchical Random Effects Model with drift and minorization conditions established in [Jones & Hobert 2004]. Even more general models are feasible in this setting, in particular in the recent paper [1] Johnson and Jones established drift and minorization conditions for a bayesian hierarchical version of a general linear mixed model.

- Establishing drift conditions might be difficult. A good first try may be  $V(x)$  proportional to  $\pi(x)^{-1/2}$  or to some suitable quadratic function.

## 5.8 Appendix - Formulas for $\rho$ and M

In the sequel the term *atomic case* and *nonatomic case* refers to  $\tilde{\beta} = 1$  and  $\tilde{\beta} < 1$  respectively. If  $\tilde{\beta} < 1$ , define

$$\alpha_1 = 1 + \frac{\log \frac{K-\tilde{\beta}}{1-\tilde{\beta}}}{\log \lambda^{-1}}, \quad \alpha_2 = \begin{cases} 1, & \text{if } \nu(C) = 1, \\ 1 + \frac{\log \tilde{K}}{\log \lambda^{-1}}, & \text{if } \nu(C) + \int_{C^c} V d\nu \leq \tilde{K}, \\ 1 + (\log \frac{K}{\tilde{\beta}}) / (\log \lambda^{-1}), & \text{otherwise.} \end{cases}$$

Then let

$$R_0 = \min\{\lambda^{-1}, (1 - \tilde{\beta})^{-1/\alpha_1}\}, \quad L(R) = \begin{cases} \frac{\tilde{\beta} R^{\alpha_2}}{1 - (1 - \tilde{\beta}) R^{\alpha_1}}, & \text{if } 1 < R < R_0, \\ \infty & \text{if } R = R_0. \end{cases}$$

### 5.8.1 Formulas for general operators

For  $\beta > 0$ ,  $R > 1$  and  $L > 1$ , let  $R_1 = R_1(\beta, R, L)$  be the unique solution  $r \in (1, R)$  of the equation

$$\frac{r-1}{r(\log(R/r))^2} = \frac{e^2 \beta (R-1)}{8(L-1)}$$

and for  $1 < r < R_1$ , define

$$K_1(r, \beta, R, L) = \frac{2\beta + 2(\log N)(\log(R/r))^{-1} - 8Ne^{-2}(r-1)r^{-1}(\log(R/r))^{-2}}{(r-1)[\beta - 8Ne^{-2}(r-1)r^{-1}(\log(R/r))^{-2}]},$$

where  $N = (L-1)/(R-1)$ .

For the *atomic case* we have  $\rho = 1/R_1(\beta, \lambda^{-1}, \lambda^{-1}K)$  and for  $\rho < \gamma < 1$ ,

$$M = \frac{\max(\lambda, K - \lambda/\gamma)}{\gamma - \lambda} + \frac{K(K - \lambda/\gamma)}{\gamma(\gamma - \lambda)} K_1(\gamma^{-1}, \beta, \lambda^{-1}, \lambda^{-1}K) + \frac{(K - \lambda/\gamma) \max(\lambda, K - \lambda)}{(\gamma - \lambda)(1 - \lambda)} + \frac{\lambda(K - 1)}{(\gamma - \lambda)(1 - \lambda)}. \quad (5.48)$$

For the *nonatomic case* let  $\tilde{R} = \arg \max_{1 < R < R_0} R_1(\beta, R, L(R))$ . Then we have  $\rho = 1/R_1(\beta, \tilde{R}, L(\tilde{R}))$  and for  $\rho < \gamma < 1$ ,

$$M = \frac{\gamma^{-\alpha_2-1}(K\gamma - \lambda)}{(\gamma - \lambda)[1 - (1 - \tilde{\beta})\gamma^{-\alpha_1}]^2} \times \left( \frac{\tilde{\beta} \max(\lambda, K - \lambda)}{1 - \lambda} + \frac{(1 - \tilde{\beta})(\gamma^{-\alpha_1} - 1)}{\gamma^{-1} - 1} \right) + \frac{\max(\lambda, K - \lambda/\gamma)}{\gamma - \lambda} + \frac{\tilde{\beta}\gamma^{-\alpha_2-2}K(K\gamma - \lambda)}{(\gamma - \lambda)[1 - (1 - \tilde{\beta})\gamma^{-\alpha_1}]^2} K_1(\gamma^{-1}, \beta, \tilde{R}, L(\tilde{R})) + \frac{\gamma^{-\alpha_2}\lambda(K - 1)}{(1 - \lambda)(\gamma - \lambda)[1 - (1 - \tilde{\beta})\gamma^{-\alpha_1}]} + \frac{K[K\gamma - \lambda - \tilde{\beta}(\gamma - \lambda)]}{\gamma^2(\gamma - \lambda)[1 - (1 - \tilde{\beta})\gamma^{-\alpha_1}]} + \frac{K - \lambda - \tilde{\beta}(1 - \lambda)}{(1 - \lambda)(1 - \gamma)} \left( (\gamma^{-\alpha_2} - 1) + (1 - \tilde{\beta})(\gamma^{-\alpha_1} - 1)/\tilde{\beta} \right). \quad (5.49)$$

## 5.8.2 Formulas for self-adjoint operators

A Markov chain is said to be reversible with respect to  $\pi$  if  $\int_{\mathcal{X}} Pf(x)g(x)\pi(dx) = \int_{\mathcal{X}} f(x)Pg(x)\pi(dx)$  for all  $f, g \in L^2(\pi)$ . For reversible Markov chains the following tighter bounds are available.

For the *atomic case* define

$$R_2 = \begin{cases} \min\{\lambda^{-1}, r_s\}, & \text{if } K > \lambda + 2\beta, \\ \lambda^{-1}, & \text{if } K \leq \lambda + 2\beta, \end{cases}$$

where  $r_s$  is the unique solution of  $1 + 2\beta r = r^{1+(\log K)(\log \lambda^{-1})}$ . Then  $\rho = R_2^{-1}$  and for  $\rho < \gamma < 1$  take  $M$  as in (5.48) with  $K_1(\gamma^{-1}, \beta, \lambda^{-1}, \lambda^{-1}K)$  replaced by  $K_2 = 1 + 1/(\gamma - \rho)$ .

For the *nonatomic case* let

$$R_2 = \begin{cases} r_s, & \text{if } L(R_0) > 1 + 2\beta R_0, \\ R_0, & \text{if } L(R_0) \leq 1 + 2\beta R_0, \end{cases}$$

where  $r_s$  is the unique solution of  $1 + 2\beta r = L(r)$ . Then  $\rho = R_2^{-1}$  and for  $\rho < \gamma < 1$  take  $M$  as in (5.49) with  $K_1(\gamma^{-1}, \beta, \tilde{R}, L(\tilde{R}))$  replaced by  $K_2 = 1 + \sqrt{\tilde{\beta}}/(\gamma - \rho)$ .

### 5.8.3 Formulas for self-adjoint positive operators

A Markov chain is said to be positive if  $\int_{\mathcal{X}} Pf(x)f(x)\pi(dx) \geq 0$  for every  $f \in L^2(\pi)$ . For reversible and positive markov chains take  $M$ 's as in Section 5.8.2 with  $\rho = \lambda$  in the *atomic case* and  $\rho = R_0^{-1}$  in the *nonatomic case*.

# Chapter 6

## Convergence Results for Adaptive Monte Carlo

Ergodicity results for adaptive Monte Carlo algorithms usually assume *time-stability* of transition kernels. On the other hand, a large class of time-inhomogeneous Markov Chains is ergodic. This suggests existence of adaptive MC algorithms which fail to satisfy the *time-stability* condition but are still ergodic. We present a modification of Atchadé-Rosenthal ergodicity Theorems (3.1 and 3.2 in [Atchadé & Rosenthal 2005]) that does not assume *time-stability* of transition kernels. We use a weaker *path-stability* condition instead, that results from *time-stability* condition by the triangle inequality.

### 6.1 Introduction

As before, we deal with computation of analytically intractable integral

$$I = \int_{\mathcal{X}} f(x)\pi(x)dx.$$

For computational efficiency of the Markov chain Monte Carlo approach, the simulated Markov chain should converge to its stationary distribution reasonably quickly. This can sometimes be achieved by careful design of the transition kernel  $P$  of the chain, on the basis of a detailed preliminary analysis of  $\pi$ . Intuitively, the more features of  $\pi$  are known, the better  $P$  can be designed. So a non-Markovian approach might be to allow the transition kernel of the simulated stochastic process  $(X_n)_{n \geq 0}$  to adapt whenever new

features of  $\pi$  are encountered during the process run. Simulations show that this approach can indeed outperform algorithms based on classical ideas. For numerous examples and an insight of how to tune the transition kernel "on the fly" see [Roberts & Rosenthal 2006] and references therein. However, since in this case  $(X_n)_{n \geq 0}$  is not a Markov chain any more, it may fail to converge to the expected asymptotic distribution even if each participating transition kernel is ergodic and has the same stationary distribution. A simple but nonintuitive example is given in Section 6.2. Difficulty to obtain general ergodicity results appears to be the main problem in adaptive Monte Carlo.

For versions of adaptive MC and related work we refer to e.g. [Fishman 1996], [Evans 1991], [Gelfand & Sahu 1994]. In more recent papers [Gilks et al. 1998] showed adaptation of the transition kernel can be performed (without damaging the ergodicity of the algorithm) on regeneration times. The idea of adaptive MC through regeneration was then investigated in [Brockwell & Kadane 2005] and [Sahu & Zhigljavsky 2003]. Convergence results in fairly general setting have been derived in [Haario et al. 2001] which was followed by refined theorems in [Atchadé & Rosenthal 2005] and a discrete state space version of those results presented in [Kohn & Nott 2005].

In each of the above mentioned papers ergodicity results either on regeneration times, or fit within the so called diminishing adaptation framework and assume the *time-stability* condition for transition kernels. Yet the existence of ergodic inhomogeneous Markov chains suggests the *time-stability* of transition kernels is not necessary for ergodicity of adaptive MC algorithms. After introductory examples in Section 6.2, in Section 6.3 we give ergodicity theorems that use a weaker *path-stability* condition, which results from the *time-stability* condition by triangle inequality. However we have to pay the price for it and formulate the *uniform ergodicity* condition in the time inhomogeneous setting, which makes it more complicated than in the original Atchadé and Rosenthal's theorems. In Section 6.4 we prove the main result of this Chapter.

## 6.2 One Intuitive and One Not-so-Intuitive Example

We begin with a simple example where we briefly analyze two stochastic processes using the same two transition matrices.

Consider the state space  $\mathcal{X} = \{0, 1\}$  and  $\pi$ , the uniform distribution on  $\mathcal{X}$ . Let

$$P_1 = \begin{bmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix} \quad \text{and} \quad P_2 = (1 - \varepsilon) \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \varepsilon P_1 \quad \text{for some } \varepsilon > 0.$$

Note that  $\pi$  is the stationary distribution for both,  $P_1$  and  $P_2$ . Let  $\varphi$  be some probability distribution on  $\{P_1, P_2\}$ . Let  $P^{(0)}, P^{(1)}, P^{(2)}, \dots$  be an iid sample from  $\varphi$ . In the sequel we will use the convention  $\min \emptyset = \infty$  and  $\max \emptyset = -\infty$ .

**Example 6.2.1.** Let  $(X_n)_{n \geq 0}$  be a stochastic process with an initial distribution  $p_0$ , evolving in step  $k$  according to the transition matrix  $P^{(k)}$ .  $(X_n)_{n \geq 0}$  is clearly an in-homogeneous Markov Chain and  $p_n$  (the distribution of  $X_n$ ) converges to the stationary distribution  $\pi$ : let  $U_n := \{k : k \leq n, P^{(k)} = P_1\}$  and  $u_n = \max U_n$ . The distribution of  $X_n$ , given  $u_n \neq -\infty$  is  $\pi$ , so we have the following bound on the total variation distance between  $p_n$  and  $\pi$ :

$$\|p_n - \pi\|_{tv} \leq P(u_n = -\infty) \xrightarrow{n \rightarrow \infty} 0 \quad a.s.$$

**Example 6.2.2.** (due to W. Niemi). Now consider  $(Y_n)_{n \geq 0}$  with an initial distribution  $q_0$  and an initial transition matrix  $Q_0$ , evolving for  $n \geq 1$  according to the following adaptive rule:

$$Q_k = \begin{cases} P_1 & \text{if } Y_{k-1} = 0 \\ P_2 & \text{if } Y_{k-1} = 1 \end{cases}$$

Note that after two consecutive 1 (and this occurs with probability at least  $\frac{1}{4}$  for any  $k, k+1$ )  $Y_n$  is trapped in 1 and can escape only with probability  $\varepsilon$ . Let  $\bar{q}_1 = \lim_{n \rightarrow \infty} P(Y_n = 1)$  and  $\bar{q}_0 = \lim_{n \rightarrow \infty} P(Y_n = 0)$ . Now it is clear, that for small  $\varepsilon$  we will have  $\bar{q}_1 \gg \bar{q}_0$  and the procedure fails to give the expected asymptotic distribution.

Both processes  $(X_n)_{n \geq 0}$  and  $(Y_n)_{n \geq 0}$  are allowed to use essentially different transition matrices in two consecutive steps. But one of them converges to the desired distribution  $\pi$  and the other one fails to converge. In our opinion it is not the "time stability" condition, that is crucial for convergence of an adaptive Monte Carlo algorithm. It is the "path-stability" condition, that reads "if the path is similar, the transition kernel should be similar as well". Obviously  $(X_n)_{n \geq 0}$  satisfies this condition and  $(Y_n)_{n \geq 0}$  does not.

In the following section we will try to formalize this intuition.

## 6.3 Convergence Results

We will similarly as in [Atchadé & Rosenthal 2005] analyze a stochastic process  $(X_n)_{n \geq 0}$  on a general state space  $\mathcal{X}$ , generated by the following algorithm:

**Algorithm 6.3.1.** *Assuming we have an initial transition kernel  $P_{x_0}$  and an initial point  $x_0 \in \mathcal{X}$ , the algorithm proceeds as follows:*

1. *If for time  $n \geq 0$  we have  $X_n = x$  and a transition kernel  $P_{n, \tilde{X}_n}$ , which is allowed to depend on the path  $\tilde{X}_n = (X_0, \dots, X_n) \in \mathcal{X}^{n+1}$ ; then sample from  $P_{n, \tilde{X}_n}(x, \cdot)$ .*
2. *Use  $\tilde{X}_{n+1} = (X_0, \dots, X_{n+1})$  to build a new transition kernel  $P_{n, \tilde{X}_{n+1}}$  to be used at time  $n + 1$ .*

For  $(X_n)_{n \geq 0}$  generated by Algorithm 6.3.1 we shall write  $P_\mu$  to denote its distribution on  $(\mathcal{X}^\infty, \mathcal{F}^\infty)$  when  $X_0 \sim \mu$ , and  $E_\mu$  to denote the expectation with respect to  $P_\mu$ . If  $\mu = \delta_x$ , we usually write  $E_x$  and  $P_x$  instead of  $E_\mu$  and  $P_\mu$ . By  $P_{\mu, n}$  we will denote the marginal distribution of  $X_n$  induced by  $P_\mu$ , thus  $P_{\mu, n}$  is a probability measure on  $\mathcal{X}$ . To denote two trajectories of length  $n + k + 1$ , that have a common initial part of length  $n + 1$  and then split, we will write  $(\tilde{x}_n, \tilde{y}_k)$  and  $(\tilde{x}_n, \tilde{y}'_k)$ .

We will prove ergodicity theorems similar to Theorem 3.1 and 3.2 in [Atchadé & Rosenthal 2005], but under modified assumptions.

**Assumption 6.3.2.** *There exist a measurable function  $V : \mathcal{X} \rightarrow [1, \infty)$  and real number sequences  $(\tau_n), (a_n), (R_n)$ , such that  $(\tau_n), (R_n) \rightarrow 0$  as  $n \rightarrow \infty$  and:*

- A.1 (uniform ergodicity) *For all  $j \geq 1, n \geq 0, x \in \mathcal{X}$  and  $\tilde{x}_n \in \mathcal{X}^{n+1}$ , there exists  $\tilde{y}'_j = (y'_1, \dots, y'_j)$  and  $0 \leq l \leq j - 1$  such that*

$$\left\| \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, y'_1, \dots, y'_i)}(x, \cdot) - \pi_{n+l, (\tilde{x}_n, y'_1, \dots, y'_l)}(\cdot) \right\|_V \leq R_j V(x). \quad (6.1)$$

- A.2 (path-stability) *For all  $x \in \mathcal{X}, \tilde{x}_n \in \mathcal{X}^{n+1}$ , there exists  $\tilde{y}'_k \in \mathcal{X}^k$ , such that  $\tilde{x}_n$  and  $\tilde{y}'_k$  satisfy (6.1) with  $j = k$  and for all  $\tilde{y}_k \in \mathcal{X}^k$ ,*

$$\left\| P_{n+k, (\tilde{x}_n, \tilde{y}_k)}(x, \cdot) - P_{n+k, (\tilde{x}_n, \tilde{y}'_k)}(x, \cdot) \right\|_V \leq K_1 \tau_n a_k V(x). \quad (6.2)$$

A.3 For all  $x \in \mathcal{X}$ ,  $\tilde{x}_n \in \mathcal{X}^{n+1}$ ,  $\tilde{y}_k \in \mathcal{X}^k$ ,

$$\left\| \pi_{n+k,(\tilde{x}_n, \tilde{y}_k)} - \pi_{n, \tilde{x}_n} \right\|_V \leq K_2 \tau_n a_k. \quad (6.3)$$

A.4 For all  $n \geq 1$ ,

$$\begin{aligned} \int V^2(x_n) P_{\mu, n}(dx_n) &= \\ &= \int \dots \int V^2(x_n) P_{n-1, \tilde{x}_{n-1}}(x_{n-1}, dx_n) \dots P_{0, \tilde{x}_0}(x_0, dx_1) \leq K_3 V^2(x_0) \end{aligned} \quad (6.4)$$

and

$$\sup_{n, \tilde{x}_n} \pi_{n, \tilde{x}_n}(V) < \infty. \quad (6.5)$$

A.5 For any finite constants  $c_1, c_2$ , define

$$B(c_1, c_2, n) := \min_{1 \leq k \leq n} (c_1 \phi_k \tau_{n-k} + c_2 R_k),$$

where  $\phi_n = \sum_{k=1}^n a_k$ . Assume that  $B(c_1, c_2, n) = \mathcal{O}(\frac{1}{n^\varepsilon})$  for some  $\varepsilon > 0$ .

Under these assumptions we will prove two ergodicity theorems:

**Theorem 6.3.3.** Let  $(X_n)_{n \geq 0}$  be the stochastic process generated by Algorithm 6.3.1 with  $X_0 = x_0$ . Under A.1-A.4 there exist constants  $k_1, k_2 < \infty$  such that for any measurable function  $f : \mathcal{X} \rightarrow R$  with  $|f| \leq V$ ,

$$\left| E_{x_0}(f(X_n) - \pi_{n, \tilde{X}_n}(f)) \right| \leq B(k_1, k_2, n) V(x_0). \quad (6.6)$$

**Theorem 6.3.4.** Under A.1-A.5, for any measurable function  $f : \mathcal{X} \rightarrow R$  and  $|f| \leq V$ , for any starting point  $x_0 \in \mathcal{X}$ ,

$$\frac{1}{n} \sum_{i=0}^{n-1} (f(X_i) - \pi_{i, \tilde{x}_i}(f)) \rightarrow 0, \quad \text{as } n \rightarrow \infty, \quad P_{x_0} \text{- a.s.} \quad (6.7)$$

*Remark 6.3.5.* 1. If  $\pi_{n, \tilde{x}_n} \equiv \pi$ , as it usually occurs in Monte Carlo setting ( $\pi$  is the invariant target distribution), then Theorem 6.3.3 gives a bound on the rate of convergence of the distribution of  $X_n$  to  $\pi$  and Theorem 6.3.4 provides a law of large numbers type result.



2. In this typical case ( $\pi_{n,\tilde{x}_n} \equiv \pi$ ) the theory of inhomogeneous Markov chains can be applied to check Assumption A.1 (compare [Douc et al. 2003]).
3. In particular this theorems can be applied in case when  $\pi_{n,\tilde{x}_n} \equiv \pi$  and  $P_{n,\tilde{x}_n} \geq \varepsilon\pi$ , for some  $\varepsilon > 0$ , as considered in [Kohn & Nott 2005].
4. Assumptions used here differ from those in [Atchadé & Rosenthal 2005], where A.1 and A.2 are as follows:

A.1' (*uniform ergodicity*) For all  $j > 0$ ,  $n \geq 0$ ,  $x \in \mathcal{X}$  and  $\tilde{x}_n \in \mathcal{X}^{n+1}$ ,

$$\|P_{n,\tilde{x}_n}^j(x, \cdot) - \pi_{n,\tilde{x}_n}(\cdot)\|_V \leq R_j V(x).$$

A.2' (*time-stability*) For all  $x \in \mathcal{X}$ ,  $\tilde{x}_n \in \mathcal{X}^{n+1}$ ,  $\tilde{y}_k \in \mathcal{X}^k$ ,

$$\|P_{n+k,(\tilde{x}_n,\tilde{y}_k)}(x, \cdot) - P_{n,\tilde{x}_n}(x, \cdot)\|_V \leq K_1 \tau_n a_k V(x).$$

The *path-stability* condition results from the *time-stability* condition by the triangle inequality, so assumption A.2 presented here is weaker. Assumptions A.1 here and A.1' in [Atchadé & Rosenthal 2005] are incomparable.  $\|P_{n,\tilde{x}_n}^j(x, \cdot) - \pi_{n,\tilde{x}_n}(\cdot)\|_V$  does not have to converge even if A.1-5 hold. It involves some computation, similar to this in the proof of Lemma 6.4.1, to show A.1', A.2' together with A.3-4 imply

$$\left\| \prod_{i=0}^{j-1} P_{n+i,(\tilde{x}_n,y'_1,\dots,y'_i)}(x, \cdot) - \pi_{n+l,(\tilde{x}_n,y'_1,\dots,y'_l)}(\cdot) \right\|_V \leq B(k_1, k_2, j)V(x),$$

so if additionally A.5 holds,

$$\left\| \prod_{i=0}^{j-1} P_{n+i,(\tilde{x}_n,y'_1,\dots,y'_i)}(x, \cdot) - \pi_{n+l,(\tilde{x}_n,y'_1,\dots,y'_l)}(\cdot) \right\|_V = \mathcal{O}\left(\frac{1}{n^\varepsilon}\right).$$

However our version is more complicated and might turn out to be difficult to check even if  $\pi_{n,\tilde{x}_n} \equiv \pi$ .

5. *Path-stability* instead of *time-stability* condition enables to apply this ergodicity theorems to Monte Carlo algorithms that are inhomogeneous in their nature, like simulated annealing. In other words we can adapt Monte Carlo methods based on inhomogeneous Markov chains as well.

6. Finally, the theorem handles our introductory toy examples i.e.  $(X_n)_{n \geq 0}$  that converges to the desired distribution satisfies A.1-A.4 (but does not satisfy A.2' in [Atchadé & Rosenthal 2005]).  $(Y_n)_{n \geq 0}$  that fails to converge, fails to satisfy assumption A.2 as well.

## 6.4 Proofs

We now proceed to prove theorems from Section 6.3. The proof follows closely Atchadé and Rosenthal [Atchadé & Rosenthal 2005]. Crucial point of the proof is Lemma 6.4.1. Once Lemma 6.4.1 is shown under our modified assumptions, we derive Theorems 6.3.3 and 6.3.4 in essentially identical manner as in [Atchadé & Rosenthal 2005]. This part of the proof is purely expository and presented here for the sake of completeness.

Let  $(\mathcal{F}_n)_{n=-\infty}^{\infty}$  be a filtration defined by:

$$\mathcal{F}_n := \begin{cases} \{\cdot, \Omega\} & \text{if } n < 0 \\ \sigma(X_0, \dots, X_n) & \text{if } n \geq 0 \end{cases} \quad (6.8)$$

and  $g_{k, \tilde{X}_k}(x) := f(x) - \pi_{k, \tilde{X}_k}(f)$ .

**Lemma 6.4.1.** *Assume A.1-A.4 hold. Then there are some constants  $0 < k_1, k_2 < \infty$  such that for any  $n \geq 0$ ,  $j \geq 1$  and any measurable function  $f$  with  $|f| \leq V$ , we have:*

$$\left\| E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j}) | \mathcal{F}_n) \right\|_2 \leq B(k_1, k_2, j)V(x_0). \quad (6.9)$$

The proof of Lemma 6.4.1 is given later in this section. We start with Theorems 6.3.3 and 6.3.4.

*Proof of Theorem 6.3.3.* Let  $n = 0$  in Lemma 6.4.1. We obtain the following:

$$\left\| E_{x_0}(g_{j, \tilde{X}_j}(X_j) | \mathcal{F}_0) \right\|_2 = |E_{x_0}(f(X_j) - \pi_{j, \tilde{X}_j}(f))| \leq B(k_1, k_2, j)V(x_0),$$

for all  $|f| \leq V$ , which is Theorem 6.3.3.  $\square$

*Proof of Theorem 6.3.4.* To prove Theorem 6.3.4 we will use the theory of mixingales. Theorem 6.5.2 used here is presented in Appendix. Let

$$Y_n := f(X_n) - \pi_{n, \tilde{X}_n}(f) - E_{x_0}(f(X_n) - \pi_{n, \tilde{X}_n}(f)). \quad (6.10)$$

The proof will proceed according to the following plan:

1. Show that

$$E_{x_0}(f(X_n) - \pi_{n, \tilde{X}_n}(f)) \rightarrow 0 \text{ as } n \rightarrow \infty. \quad (6.11)$$

2. Show that  $(Y_n)_{n \geq 0}$  is a mixingale of size  $-\frac{\varepsilon}{2}$  and use Theorem 6.5.2 to conclude that

$$\frac{1}{n} \sum_{i=0}^{n-1} Y_i \rightarrow 0 \text{ as } n \rightarrow \infty \quad P_{x_0} \text{ a.s.} \quad (6.12)$$

3. The foregoing results in

$$\frac{1}{n} \sum_{i=0}^{n-1} (f(X_i) - \pi_{i, \tilde{X}_i}(f)) \rightarrow 0 \text{ as } n \rightarrow \infty \quad P_{x_0} \text{ a.s.} \quad (6.13)$$

This states Theorem 6.3.4.

To see that (6.11) holds, it is enough to recall Theorem 6.3.3 and Assumption A.5.

To prove (6.12) consider first condition (6.30). Since the filtration is defined by (6.8), we have  $E(Y_n | \mathcal{F}_{n+j}) = Y_n$  and (6.30) is satisfied for any positive number sequences  $(c_n)$  and  $(\psi_n)$ .

Condition (6.29) is obviously satisfied for  $j \geq n$ , for any positive number sequences  $(c_n)$  and  $(\psi_n)$  as well, since  $EY_n = 0$ . For the case  $j < n$  we will use Lemma 6.4.1:

$$\begin{aligned} \|E_{x_0}(Y_n | \mathcal{F}_{n-j})\|_2 &= \left\| E_{x_0} \left( g_{n, \tilde{X}_n}(X_n) - (E_{x_0}(g_{n, \tilde{X}_n}(X_n))) | \mathcal{F}_{n-j} \right) \right\|_2 \\ &\leq \left\| E_{x_0} \left( g_{n, \tilde{X}_n}(X_n) | \mathcal{F}_{n-j} \right) \right\|_2 + \\ &\quad + \left\| E_{x_0} \left( E_{x_0}(g_{n, \tilde{X}_n}(X_n)) | \mathcal{F}_{n-j} \right) \right\|_2 \\ &= \left\| E_{x_0} \left( g_{n, \tilde{X}_n}(X_n) | \mathcal{F}_{n-j} \right) \right\|_2 + \left\| E_{x_0} \left( g_{n, \tilde{X}_n}(X_n) | \mathcal{F}_0 \right) \right\|_2 \\ &\leq B(k_1, k_2, j)V(x_0) + B(k_1, k_2, n)V(x_0) \\ &= \mathcal{O}(j^{-\varepsilon}) + \mathcal{O}(n^{-\varepsilon}) = \mathcal{O}(j^{-\varepsilon}) \end{aligned}$$

Now we set in (6.29)  $c_n \equiv 1$  and take appropriate  $\psi_j$ , such that  $\psi_j = \mathcal{O}(j^{-\varepsilon})$ . Hence  $(Y_n)_{n \geq 0}$  is a mixingale of size  $-\frac{\varepsilon}{2}$ . Since  $\frac{c_n}{n} = \mathcal{O}(n^{-1})$  and  $-1 < \min\{-\frac{1}{2}, \frac{\varepsilon}{2} - 1\}$ , we can apply Theorem 6.5.2 and conclude that  $\frac{1}{n} \sum_{i=0}^{n-1} Y_i \rightarrow 0$  as  $n \rightarrow \infty$   $P_{x_0}$  a.s.

Combining (6.11) and (6.12), we get (6.13) by an elementary argument.  $\square$

Now we proceed to prove Lemma 6.4.1.

*Proof of Lemma 6.4.1.* Note that

$$\pi_{n, \tilde{X}_n}(g_{n, \tilde{X}_n}) = \pi_{n, \tilde{X}_n}(f - \pi_{n, \tilde{X}_n}(f)) = 0 \quad P_{x_0} \text{ a.s.} \quad (6.14)$$

The idea of the proof is to split the quantity  $\|E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_n)\|_2$  into two terms, say  $A$  and  $B$  and bound them using Assumptions A.1, A.2 and A.3.

Denote by  $(\tilde{x}_n, \tilde{y}_j) = (\tilde{x}_n, y_1, \dots, y_j)$  a trajectory of length  $n+j$ . According to this notation we will usually write  $y_i$  for  $x_{n+i}$ . Given  $(X_0, \dots, X_n) = \tilde{x}_n$  we have

$$\begin{aligned} E_{x_0}(g_{n, \tilde{X}_n}(X_{n+j})|\tilde{X}_n = \tilde{x}_n) &= \\ &= \int g_{n, \tilde{x}_n}(y_j) P_{n+j-1, (\tilde{x}_n, y_1, \dots, y_{j-1})}(y_{j-1}, dy_j) \dots P_{n, \tilde{x}_n}(x_n, dy_1) \\ &= \eta_{j-1}(\tilde{x}_n) + \\ &+ \int g_{n, \tilde{x}_n}(y_j) P_{n+j-1, (\tilde{x}_n, \tilde{y}'_{j-1})}(y_{j-1}, dy_j) P_{n+j-2, (\tilde{x}_n, \tilde{y}_{j-2})}(y_{j-2}, dy_{j-1}) \dots \\ &\dots P_{n, \tilde{x}_n}(x_n, dy_1), \end{aligned}$$

where  $\tilde{y}'_j = (y'_1, \dots, y'_j)$  is as in Assumption A.2, and

$$\begin{aligned} \eta_{j-1}(\tilde{x}_n) &= \\ &= \int g_{n, \tilde{x}_n}(y_j) \left( P_{n+j-1, (\tilde{x}_n, \tilde{y}_{j-1})}(y_{j-1}, dy_j) - P_{n+j-1, (\tilde{x}_n, \tilde{y}'_{j-1})}(y_{j-1}, dy_j) \right) \\ &P_{n+j-2, (\tilde{x}_n, \tilde{y}_{j-2})}(y_{j-2}, dy_{j-1}) \dots P_{n, \tilde{x}_n}(x_n, dy_1). \end{aligned}$$

By exchanging transition kernels for all coordinates, we get:

$$E_{x_0} \left( g_{n, \tilde{X}_n} (X_{n+j}) | \tilde{X}_n = \tilde{x}_n \right) = \sum_{k=1}^{j-1} \eta_k(\tilde{x}_n) + \left( \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)} \right) g_{n, \tilde{x}_n}(x_n), \quad (6.15)$$

where

$$\begin{aligned} \eta_k(\tilde{x}_n) &= \int \left( \prod_{i=k+1}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)} \right) g_{n, \tilde{x}_n}(y_{k+1}) \\ &\quad \left( P_{n+k, (\tilde{x}_n, \tilde{y}'_k)}(y_k, dy_{k+1}) - P_{n+k, (\tilde{x}_n, \tilde{y}'_k)}(y_k, dy_{k+1}) \right) \\ &\quad P_{n+k-1, (\tilde{x}_n, \tilde{y}'_{k-1})}(y_{k-1}, dy_k) \dots \dots P_{n, \tilde{x}_n}(x_n, dy_1). \end{aligned} \quad (6.16)$$

Consider the second term of the right hand side of (6.15):

$$\begin{aligned} &\left| \left( \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)} \right) g_{n, \tilde{x}_n}(x_n) \right| = \\ &= \left| \left( \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)} \right) (f - \pi_{n, \tilde{x}_n}(f)) \right| \\ &= \left| \left( \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)} \right) f(x_n) - \pi_{n, \tilde{x}_n}(f) \right| \\ &= \left| \left( \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)} \right) f(x_n) - \pi_{n+l, (\tilde{x}_n, \tilde{y}'_l)}(f) \right| + \left| \pi_{n+l, (\tilde{x}_n, \tilde{y}'_l)}(f) - \pi_{n, \tilde{x}_n}(f) \right| \\ &\leq \left\| \prod_{i=0}^{j-1} P_{n+i, (\tilde{x}_n, \tilde{y}'_i)}(x_n, \cdot) - \pi_{n+l, (\tilde{x}_n, \tilde{y}'_l)}(\cdot) \right\|_V + \left\| \pi_{n+l, (\tilde{x}_n, \tilde{y}'_l)} - \pi_{n, \tilde{x}_n} \right\|_V \\ &\leq R_j V(x_n) + K_2 \tau_n a_j, \end{aligned} \quad (6.17)$$

where the inequalities result from Assumptions A.1 and A.3.

We will now bound the first term of the right hand side of (6.15). Note that since  $g_{n, \tilde{x}_n}(y_{k+1}) = f(y_{k+1}) - \pi_{n, \tilde{x}_n}(f)$  and  $\pi_{n, \tilde{x}_n}(f)$  given  $\tilde{x}_n$  is some real

number, we obtain:

$$\begin{aligned} \left( \prod_{i=k+1}^{j-1} P_{n+i,(\tilde{x}_n, \tilde{y}'_i)} \right) g_{n, \tilde{x}_n}(y_{k+1}) &= \\ &= \left( \prod_{i=k+1}^{j-1} P_{n+i,(\tilde{x}_n, \tilde{y}'_i)} \right) f(y_{k+1}) - \pi_{n, \tilde{x}_n}(f) \end{aligned} \quad (6.18)$$

and

$$\begin{aligned} \int \pi_{n, \tilde{x}_n}(f) \left( P_{n+k,(\tilde{x}_n, \tilde{y}_k)}(y_k, dy_{k+1}) - P_{n+k,(\tilde{x}_n, \tilde{y}'_k)}(y_k, dy_{k+1}) \right) \\ P_{n+k-1,(\tilde{x}_n, \tilde{y}_{k-1})}(y_{k-1}, dy_k) \dots \dots P_{n, \tilde{x}_n}(x_n, dy_1) = 0. \end{aligned} \quad (6.19)$$

Hence using (6.18) and (6.19) we get

$$\begin{aligned} \eta_k(\tilde{x}_n) &= \int \left( \prod_{i=k+1}^{j-1} P_{n+i,(\tilde{x}_n, \tilde{y}'_i)} \right) f(y_{k+1}) \\ &\quad \left( P_{n+k,(\tilde{x}_n, \tilde{y}_k)}(y_k, dy_{k+1}) - P_{n+k,(\tilde{x}_n, \tilde{y}'_k)}(y_k, dy_{k+1}) \right) \\ &\quad P_{n+k-1,(\tilde{x}_n, \tilde{y}_{k-1})}(y_{k-1}, dy_k) \dots \dots P_{n, \tilde{x}_n}(x_n, dy_1). \end{aligned} \quad (6.20)$$

Since for each  $0 \leq l \leq j - k - 2$  and  $\tilde{y}''_l$  we have

$$\begin{aligned} \left| \left( \prod_{i=k+1}^{j-1} P_{n+i,(\tilde{x}_n, \tilde{y}'_i)} \right) f(y_{k+1}) \right| &= \left| \left( \prod_{i=k+1}^{j-1} P_{n+i,(\tilde{x}_n, \tilde{y}'_i)} \right) f(y_{k+1}) \right. \\ &\quad \left. - \pi_{n+k+1+l,(\tilde{x}_n, \tilde{y}'_{k+1}, \tilde{y}''_l)}(f) + \pi_{n+k+1+l,(\tilde{x}_n, \tilde{y}'_{k+1}, \tilde{y}''_l)}(f) \right| \\ &\leq \left\| \prod_{i=k+1}^{j-1} P_{n+i,(\tilde{x}_n, \tilde{y}'_i)}(y_{k+1}, \cdot) - \pi_{n+k+1+l,(\tilde{x}_n, \tilde{y}'_{k+1}, \tilde{y}''_l)}(\cdot) \right\|_V \\ &\quad + \left| \pi_{n+k+1+l,(\tilde{x}_n, \tilde{y}'_{k+1}, \tilde{y}''_l)}(f) \right|, \end{aligned} \quad (6.21)$$

we can apply A.1 and write an analogous equality to (6.19) for  $\pi_{n+k+1+l,(\tilde{x}_n, \tilde{y}'_{k+1}, \tilde{y}''_l)}(f)$

resulting from A.1 to get:

$$\begin{aligned}
|\eta_k(\tilde{x}_n)| &\leq \sup_{I:\mathcal{X}\rightarrow\{-1,1\}} \int R_{j-1-k} V(y_{k+1}) I(y_{k+1}) \\
&\quad \left( P_{n+k,(\tilde{x}_n,\tilde{y}_k)}(y_k, dy_{k+1}) - P_{n+k,(\tilde{x}_n,\tilde{y}'_k)}(y_k, dy_{k+1}) \right) \\
&\quad P_{n+k-1,(\tilde{x}_n,\tilde{y}_{k-1})}(y_{k-1}, dy_k) \dots\dots\dots P_{n,\tilde{x}_n}(x_n, dy_1) \\
&\leq R_{j-1-k} \int K_1 \tau_n a_k V(y_k) \\
&\quad P_{n+k-1,(\tilde{x}_n,\tilde{y}_{k-1})}(y_{k-1}, dy_k) \dots\dots\dots P_{n,\tilde{x}_n}(x_n, dy_1) \\
&\leq r_0 \tau_n a_k E_{x_0}(V(X_{n+k}) | \tilde{X}_n = \tilde{x}_n). \tag{6.22}
\end{aligned}$$

Where the second inequality results from *path-stability* condition A.2 and  $r_0$  is some finite constant, since  $K_1 < \infty$  and  $(R_n) \rightarrow 0$  as  $n \rightarrow \infty$ .

Putting (6.17) and (6.22) together in (6.15), we get:

$$\begin{aligned}
|E_{x_0}(g_{n,\tilde{X}_n}(X_{n+j}) | \mathcal{F}_n)| &\leq \\
&\leq R_j V(X_n) + K_2 \tau_n a_j + r_0 \tau_n \sum_{k=1}^{j-1} a_k E_{x_0}(V(X_{n+k}) | \mathcal{F}_n). \tag{6.23}
\end{aligned}$$

By Assumption A.3 we have

$$\begin{aligned}
|E_{x_0}(g_{n+j,\tilde{X}_{n+j}}(X_{n+j}) | \mathcal{F}_n)| &\leq |E_{x_0}(g_{n,\tilde{X}_n}(X_{n+j}) | \mathcal{F}_n)| \\
&\quad + E_{x_0}\left(|\pi_{n+j,\tilde{X}_{n+j}}(f) - \pi_{n,\tilde{X}_n}(f)| | \mathcal{F}_n\right) \\
&\leq |E_{x_0}(g_{n,\tilde{X}_n}(X_{n+j}) | \mathcal{F}_n)| + K_2 \tau_n a_j. \tag{6.24}
\end{aligned}$$

We now combine (6.23) and (6.24) to obtain the first inequality of the following bound:

$$\begin{aligned}
\|E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_n)\|_2 &\leq R_j \|V(X_n)\|_2 + 2K_2 \tau_n a_j + \\
&\quad + r_0 \tau_n \sum_{k=1}^{j-1} a_k \|E_{x_0}(V(X_{n+k})|\mathcal{F}_n)\|_2 \\
&\leq R_j \|V(X_n)\|_2 + \\
&\quad + \max\{r_0, 2K_2\} \tau_n \sum_{k=1}^j a_k \|V(X_{n+k})\|_2 \\
&\leq R_j \sqrt{K_3} V(X_0) + \\
&\quad + \max\{r_0, 2K_2\} \tau_n \sum_{k=1}^j a_k \sqrt{K_3} V(X_0) \\
&\leq V(x_0)(r_3 R_j + r_2 \tau_n \phi_j), \tag{6.25}
\end{aligned}$$

where we use Assumption A.4 and apply

$$\begin{aligned}
\|E_{x_0}(V(X_{n+k})|\mathcal{F}_n)\|_2 &= \{E[(E_{x_0}(V(X_{n+k})|\mathcal{F}_n))^2]\}^{1/2} \\
&\leq \{E(E_{x_0}(V^2(X_{n+k})|\mathcal{F}_n))\}^{1/2} \\
&= \{EV^2(X_{n+k})\}^{1/2} = \|V(X_{n+k})\|_2
\end{aligned}$$

The constants in (6.25) are defined as  $r_3 := \sqrt{K_3}$ ,  $r_2 := \max\{r_0, 2K_2\} \sqrt{K_3}$  and  $\phi_j := \sum_{k=1}^j a_k$ .

Since  $(\mathcal{F}_n)_{n=-\infty}^{\infty}$  is a filtration,  $\mathcal{F}_n \subseteq \mathcal{F}_{n+j-k}$ , for  $k = 1, \dots, j$  and therefore

$$E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_n) = E_{x_0}(E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_{n+j-k})|\mathcal{F}_n).$$

This implies

$$\left\{E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_n)\right\}^2 \leq E_{x_0}\left(\left\{E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_{n+j-k})\right\}^2|\mathcal{F}_n\right).$$

And therefore

$$\left\|E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_n)\right\|_2 \leq \left\|E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_{n+j-k})\right\|_2. \tag{6.26}$$

We now apply (6.25) to the right hand side of (6.26) and get:

$$\left\|E_{x_0}(g_{n+j, \tilde{X}_{n+j}}(X_{n+j})|\mathcal{F}_n)\right\|_2 \leq V(x_0)(r_3 R_k + r_2 \tau_{n+j-k} \phi_k). \tag{6.27}$$



Finally, since (6.27) holds for every  $k = 1, \dots, j$ , we can take the minimum:

$$\left\| E_{x_0}(g_{n+j, \tilde{x}_{n+j}}(X_{n+j}) | \mathcal{F}_n) \right\|_2 \leq V(x_0) \min_{1 \leq k \leq j} \{r_3 R_k + r_2 \tau_{n+j-k} \phi_k\}. \quad (6.28)$$

Obviously  $V(x_0) \min_{1 \leq k \leq j} \{r_3 R_k + r_2 \tau_{n+j-k} \phi_k\} \leq V(x_0) B(k_1, k_2, j)$  for some constants  $k_1$  and  $k_2$ , which completes the proof of the lemma.  $\square$

Hence the proof of Theorems (6.3.3) and (6.3.4) is complete as well.

## 6.5 Appendix - Mixingales

We present here a version of Strong Law of Large Numbers for mixingales that is used to conclude the proof of Theorem 6.3.4. Theorem 6.5.2 presented here is a version of Corollary 2.1 in [Davidson & Jong 1997]. For an introduction to mixingales see the books [Hall & Heyde 1980] or [Davidson 1994].

Let  $(Z_n)_{n \geq 0}$  be a real-valued stochastic process on some probability space  $(\Omega, \mathcal{F}, P)$ . Assume  $(Z_n)$  is  $L_2$ -bounded, i.e.  $\|Z_n\|_2 = \left\{ \int Z_n^2(\omega) dP(\omega) \right\}^{1/2} < \infty$  for all  $n \geq 0$ . Let  $(\mathcal{F}_n)_{n=-\infty}^{\infty}$  be a filtration.

**Definition 6.5.1.** The process  $(Z_n)_{n \geq 0}$  is a  $L^2$ -mixingale with respect to filtration  $(\mathcal{F}_n)_{n=-\infty}^{\infty}$  if there exist real number sequences  $(c_n)$  and  $(\psi_n)$ ,  $\psi_n \rightarrow 0$  as  $j \rightarrow \infty$ , such that for all  $n \geq 0$  and all  $j \geq 0$ ,

$$\left\| E(Z_n | \mathcal{F}_{n-j}) \right\|_2 \leq c_n \psi_j, \quad (6.29)$$

and

$$\left\| Z_n - E(Z_n | \mathcal{F}_{n+j}) \right\|_2 \leq c_n \psi_{j+1}. \quad (6.30)$$

If for some  $\lambda > 0$ ,  $\psi_n = \mathcal{O}(n^{-\lambda-\varepsilon})$  for some  $\varepsilon > 0$ , we say that mixingale  $Z_n$  is of size  $-\lambda$ .

**Theorem 6.5.2.** Let  $(Z_n)$  be a  $L^2$ -mixingale of size  $-\lambda$ . If  $\frac{c_n}{n} = \mathcal{O}(n^\alpha)$ , where  $\alpha < \min\{-\frac{1}{2}, \lambda - 1\}$ , then  $\frac{1}{n} \sum_{i=0}^{n-1} Z_i \rightarrow 0$  a.s.

# Bibliography

- [Aldous 1987] Aldous D., 1987, *On the Markov Chain Simulation Method for Uniform Combinatorial Distributions and Simulated Annealing*. Probability in the Engineering and Informational Sciences 1, 33-46.
- [Atchadé & Rosenthal 2005] Atchadé Y. F., Rosenthal J. S., 2005. *On Adaptive Markov Chain Monte Carlo Algorithms*. Bernoulli 11, 815–828.
- [Athreya & Ney 1978] Athreya K. B. and Ney P. (1978). A new approach to the limit theory of recurrent Markov chains. *Trans. Amer. Math. Soc.* **245** 493–501.
- [Baxendale 2005] Baxendale P. H., 2005. *Renewal Theory and Computable Convergence Rates for Geometrically Ergodic Markov Chains*. Ann. Appl. Prob. 15, 700–738.
- [Bednorz 2008] Bednorz, W., (2008) *Sharp Kendall Theorem and its Applications to Convergence Rates for Geometrically Ergodic Markov Chains*. Preprint.
- [Bednorz, Latała & Łatuszyński 2008] Bednorz W., Latała R., Łatuszyński K. (2008). *A Regeneration Proof of the Central Limit Theorem for Uniformly Ergodic Markov Chains*. Electronic Communications in Probability, 13, 85–98.
- [Bednorz & Łatuszyński 2007] Bednorz, W., Łatuszyński, K., (2007), *A few Remarks on "Fixed-Width Output Analysis for Markov Chain Monte Carlo" by Jones et al.* Journal of the American Statistical Association 102 (480), 1485-1486.
- [Billingsley 1968] Billingsley P., 1968. *Convergence of Probability Measures*. Wiley, New York.

- [Bradley 1983] Bradley, R. C. (1983) Information regularity and the central limit question. *Rocky Mountain Journal of Mathematics* **13** 77–97.
- [Breyer & Roberts 2001] Breyer L. A. and Roberts G. O. (2001). Catalytic perfect simulation. *Methodol. Comput. Appl. Probab.* **3** 161–177.
- [Brockwell & Kadane 2005] Brockwell A. E., Kadane J.B., 2005. *Identification of Regeneration Times in MCMC Simulation, with Application to Adaptive Schemes*. *Journal of Computational and Graphical Statistics*, **14**, 436–458.
- [Casella & Robert 1999] Casella G., Robert C. P., 1999. *Monte Carlo Statistical Methods*. Springer-Verlag, New York.
- [Chan & Yue 1996] Chan K. S., Yue H. (1996), "Asymptotic Efficiency of the Sample Mean in Markov Chain Monte Carlo Schemes," *Journal of the Royal Statistical Society, Series B.* **58** (3), 525-539.
- [Chow & Robbins 1965] Chow Y. S., Robbins H., (1965) *On the Asymptotic Theory of Fixed-Width Sequential Confidence Intervals for the Mean*. *The Annals of Mathematical Statistics*, **36**, 457–462.
- [Cogburn 1972] Cogburn, R. (1972). The Central Limit Theorem for Markov Processes. *In Le Cam, L. E., Neyman, J. & Scott, E. L. (Eds) Proc. Sixth Ann. Berkley Symp. Math. Sttist. and Prob.* **2** 458–512.
- [Davidson 1994] Davidson J., 1994. *Stochastic Limit Theory*. Oxford University Press, Oxford, New York.
- [Davidson & Jong 1997] Davidson J., de Jong R., 1997. *Strong laws of large numbers for dependent heterogenous processes: a synthesis of recent and new results*. *Econometric Reviews* **16**, 251-279.
- [Douc et al. 2003] Douc R., Moulines E., Rosenthal J. S., 2003. *Quantitative bounds on convergence of time-inhomogeneous Markov Chains*. *Ann. Appl. Prob.* **14**, 1643-1665.
- [Doukhan et al. 1994] Doukhan P., Massart P., Rio E., 1994. *The Functional Central Limit Theorem for Strongly Mixing Processes*. *Annales de l'Institut Henri Poincare, Section B, Calcul de Probabilites et Statistique*, **30**, 63–82.

- [Evans 1991] Evans M., 1991. *Chaining via annealing*. Ann. Statistics 19, 382-393.
- [Fishman 1996] Fishman G. S., 1996. *Monte Carlo. Concepts, algorithms and applications*. Springer.
- [Gelfand & Sahu 1994] Gelfand A. E., Sahu S. K., 1994. *On Markov chain Monte Carlo acceleration*. J. Computational and Graphical Stat. 3, 261-276.
- [Gelfand & Smith 1990] Gelfand, A. E., Smith, A. F. M., 1990, *Sampling-based Approaches to Calculating Marginal Densities*. J. Amer. Statist. Assoc. 85, 398-409.
- [Geyer 1992] Geyer C. J., 1992, *Practical Markov Chain Monte Carlo*. Stat. Sci. 7 (4), 473-511.
- [Gilks et al. 1998] Gilks W. R., Roberts G. O., Sahu S. K., 1998. *Adaptive Markov chain Monte Carlo through regeneration*. J. Amer. Statist. Assoc. 93 (443), 1045-1054.
- [Gillman 1998] Gillman D., 1998, *A Chernoff Bound for Random Walks on Expander Graphs*. SIAM J. Comput. 27 (4), 1203-1220.
- [Glynn & Ormoneit 2002] Glynn P. W., Ormoneit D. 2002 *Hoeffding's Inequality for Uniformly Ergodic Markov Chains*. Statist. and Probab. Lett. 56, 143-146.
- [Glynn & Whitt 1992] Glynn P. W., Whitt W., (1992) *The Asymptotic Validity of Sequential Stopping Rules for Stochastic Simulations*. The Annals of Applied Probability, 2, 180-198.
- [Haario et al. 2001] Haario H., Saksman E., Tamminen J., 2001. *An adaptive Metropolis algorithm*. Bernoulli 7, 223-242.
- [Hall & Heyde 1980] Hall P., Heyde C. C., 1980. *Martingale limit theory and its applications*. Academic Press, New York.
- [Häggström 2005] Häggström, O. (2005), "On the Central Limit Theorem for Geometrically Ergodic Markov Chains," *Probability Theory and Related Fields*, 132, 74-82.

- [Hobert & Geyer 1998] Hobert J.P., Geyer, C.J., 1998, *Geometric Ergodicity of Gibbs and block Gibbs samplers for Hierarchical Random Effects Model*. J. Multivariate Anal. 67, 414-430.
- [Hobert et al. 2002] Hobert, J. P., Jones G. J., Pressnell, B., Rosenthal, J.S. (2002), "On the Applicability of Regenerative Simulation in Markov Chain Monte Carlo," *Biometrika*, 89, 731-743.
- [Hobert & Robert 2004] Hobert J. P. and Robert C. P. (2004). A mixture representation of  $\pi$  with applications in Markov chain Monte Carlo and perfect smpling. *Ann. Appl. Probab.* **14** 1295–1305.
- [Ibragimov & Linnik 1971] Ibragimov, I. A. and Linnik, Y. V. (1971). *Independent and Stationary Sequences of Random Variables*. Wolters-Noordhof, Groningen.
- [1] Johnson A. A., Jones G. L., (2007) *Gibbs Sampling for a Bayesian Hierarchical Version of the General Linear Mixed Model*. Preprint.
- [Jones 2005] Jones, G. L. (2005). On the Markov chain central limit theorem. *Probability Surveys* **1** 299–320.
- [Jones et al. 2006] Jones, G. L., Haran, M., Caffo, B. S., Neath, R. (2006), "Fixed-Width Output Analysis for Markov Chain Monte Carlo," *Journal of the American Statistical Association*, 101, 1537-1547.
- [Jones & Hobert 2004] Jones G.L., Hobert J. P., 2004, *Sufficient Burn-in for Gibbs Samplers for a Hierarchical Random Effects Model*. The Annals of Statistics 32 (2), 784-817.
- [Kipnis & Varadhan 1986] Kipnis C., Varadhan S. R. S., 1986 *Central Limit Theorem for Additive Functionals of Reversible Markov Processe and Applications to Simple Exclusions* Commun. Math. Phys. 104, 1-19.
- [Kohn & Nott 2005] Kohn R., Nott D., 2005. *Adaptive Sampling for Bayesian Variable Selection*. *Biometrika*, 92, 747–763.
- [Kontoyiannis at al. 2005] Kontoyiannis I., Lastras-Montano L., Meyn S. P. 2005 *Relative Entropy and Exponential Deviation Bounds for General Markov Chains*. 2005 IEEE International Symposium on Information Theory.

- [León & Perron 2004] León C. A., Perron F., 2004. *Optimal Chernoff Bounds for Finite Reversible Markov Chains*. Ann. Appl. Prob. 14, 958-970.
- [Liu, JS 2001] Liu J. S., 2001. *Monte Carlo Strategies in Scientific Computing*. Springer.
- [Liu, W 1997] Liu W., (1997) *Improving the Fully Sequential Sampling Scheme of Anscombe-Chow-Robbins*. The Annals of Statistics, 25, 2164–2171.
- [Mengersen & Tweedie 1996] Mengersen K. L., Tweedie R. L., 1996. *Rates of Convergence of the Hastings and Metropolis Algorithms*. Annals of Statistics 24, 101-121.
- [Metropolis et al. 1953] Metropolis N., Rosenbluth A. W., Rosenbluth M. N., Teller A. H., Teller E., 1953. *Equations of state calculations by fast computing machines*. J. Chem. Phys. 21, 1087-1091.
- [Meyn & Tweedie 1993] Meyn S. P., Tweedie R. L., 1993. *Markov Chains and Stochastic Stability*. Springer-Verlag.
- [Mykland et al. 1995] Mykland P., Tierney L., Yu B., (1995) *Regeneration in Markov Chain Samplers*. Journal of the American Statistical Association, 90, 233–241.
- [Nadas 1969] Nadas A., (1969) *An Extension of a Theorem of Chow and Robbins on Sequential Confidence Intervals for the Mean*. The Annals of Mathematical Statistics, 40, 667–671.
- [Niemiro & Pokarowski 2007] Niemiro W., Pokarowski P., 2007. *Fixed Precision MCMC Estimation by Median of Products of Averages*. Submitted.
- [Nummelin 1978] Nummelin E. (1978). A splitting technique for Harris recurrent chains. *Z. Wahrscheinlichkeitstheorie und Verw. Geb.* **43** 309–318.
- [Nummelin 1984] Nummelin E. (1984). *General Irreducible Markov Chains and Nonnegative Operators*. Cambridge University Press, Cambridge.
- [Nummelin 2002] Nummelin E. (2002). MC's for MCMC'ists. *International Statistical Review.* **70** 215–240.

- [Propp & Wilson 1996] Propp, J. G. and Wilson, D. B. (1996). Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures and Algorithms*. **9** 223–252.
- [Robert 1994] Robert C. P., 1994 *The Bayesian Choice*. Springer-Verlag.
- [Roberts & Rosenthal 1997a] Roberts G. O., Rosenthal J. S., 1997. *Shift-coupling and convergence rates of ergodic averages*. Comm. in Stat. - Stoch. Models 13, 147-165.
- [Roberts & Rosenthal 1997b] Roberts G. O., Rosenthal J. S., 1997. *Geometric Ergodicity and Hybrid Markov Chains*. Elec. Comm. Prob. 2.
- [Roberts & Rosenthal 2005] Roberts G. O., Rosenthal J. S., 2005. *General state space Markov chains and MCMC algorithms*. Probability Surveys 1:20-71.
- [Roberts & Rosenthal 2006] Roberts G.O., Rosenthal J.S., *Examples of Adaptive MCMC*, Preprint (2006).
- [Roberts & Tweedie 1999] Roberts., G. O., Tweedie, R. L., 1999, *Bounds on Regeneration Times and Convergence Rates for Markov Chains*. Stochastic Process. Appl. 91, 337-338.
- [Rosenthal 1995a] Rosenthal, J. S., 1995, *Rates of Convergence for Gibbs Sampling for Variance Component Models*. The Annals of Statistics, 23, 740-761.
- [Rosenthal 1995b] Rosenthal, J. S., 1995, *Minorization Conditions and Convergence Rates for Markov Chain Monte Carlo*. Journal of the American Statistical Association, 90, 558-566.
- [Sahu & Zhigljavsky 2003] Sahu S. K., Zhigljavsky A. A., 2003. *Self-regenerative Markov Chain Monte Carlo with adaptation*. Bernoulli 9, 395-422.