

The Wang-Landau algorithm in general state spaces: Applications and convergence analysis

Yves F. Atchadé* and Jun S. Liu†

(First version Nov. 2004; Revised Feb. 2007, Aug. 2008)

Abstract: The Wang-Landau algorithm ([21]) is a recent Monte Carlo method that has generated much interest in the Physics literature due to some spectacular simulation performances. The objective of this paper is two-fold. First, we show that the algorithm can be naturally extended to more general state spaces and used to improve on Markov Chain Monte Carlo schemes of more interest in Statistics. In a second part, we study asymptotic behaviors of the algorithm. We show that with an appropriate choice of the step-size, the algorithm is consistent and a strong law of large numbers holds under some fairly mild conditions. We have also shown by simulations the potential advantage of the WL algorithm for problems in the Bayesian inference.

AMS 2000 subject classifications: Primary 60C05, 60J27, 60J35, 65C40.

Keywords and phrases: Monte Carlo methods, Wang-Landau algorithm, Multicanonical sampling, Trans-dimensional MCMC, Adaptive MCMC, Geometric ergodicity, Stochastic approximation.

1. Introduction

Although the idea of Monte Carlo computation has been around for more than a century, its first real scientific use occurred during the World War II when the first generation computer became available. Nick Metropolis coined the name “Monte Carlo” for the method when he was at Los Alamos National Labs and it quick evolved to an active research area due to the active involvements of leading physicists in the Labs. Ever since then, physicists have been at the forefront of the methodological research in the field. One of their latest additions is an algorithm proposed by F. Wang and D. P. Landau ([21]). The Wang-Landau (WL) algorithm has been successfully applied to some complex sampling problems in physics. The algorithm is closely related to *multicanonical sampling*, a method due to B. A. Berg and T. Neuhaus ([9]). Briefly,

*Department of Statistics, University of Michigan, email: yvesa@umich.edu

†Department of Statistics, Harvard University, email: jliu@stat.harvard.edu

if π is the probability measure of interest, the idea behind multicanonical sampling is to obtain an importance sampling distribution by partitioning the state space along the energy function $(-\log \pi(x))$ and re-weighting appropriately each component of the partition so that the modified distribution π^* spends equal amount of time in each component, i.e., uniform in the energy space. The method is often criticized for the difficulty involved in computing the weights. The main contribution of the WL algorithm is in proposing an efficient algorithm that *simultaneously* computes the balancing weights and samples from the re-weighted distribution.

The objective of this paper is to take a more probabilistic look at the WL algorithm and to explore its potential for Monte Carlo simulation problems of more direct interest to statisticians. We achieve this goal by proposing a general state space version of the algorithm. Then we show that the WL algorithm offers an effective strategy to improve on simulated tempering and trans-dimensional MCMC.

From a probabilistic standpoint, the WL algorithm is an interesting example of adaptive Markov Chain Monte Carlo (MCMC). Adaptive MCMC is an approach to Monte Carlo simulation where the transition kernel of the algorithm is sequentially adjusted over time in order to achieve some prescribed optimality. Some early work on the subject includes [13], [4], [15] See also [10], [19]. Early theoretical analysis includes ([6], [2], [1], [20]). We take a similar path-wise approach to analyze the WL algorithm. The analysis of the WL algorithm is not a straightforward application of the theory in the aforementioned papers because of the specific adaptive control involved. The key point is the stability of the algorithm. We say that the WL algorithm is *stable* if no component of the partition receives infinitely more visits than any other component as $n \rightarrow \infty$. On a stable sample path and under appropriate conditions, we show that the WL algorithm learns the optimal weights and satisfies a strong law of large numbers (Theorem 4.1). In the specific cases of multicanonical sampling and simulated tempering, which includes the original the WL algorithm, we show that the algorithm is stable and that the aforementioned limit results hold.

It came to our attention after the first draft of this paper that a similar extension of the WL algorithm has been proposed independently by F. Liang and coworkers (see e.g. [16]). Their approach differs from the WL approach in that these authors took a more classic approach based on stochastic approximation with step-sizes set deterministically.

Our proposed generalization to the WL algorithm is presented in Section 2. Some particular

cases are discussed in Section 3. The theoretical analysis is discussed in Section 4 but the proofs are postponed to Section 6 to facilitate the flow of ideas.

2. The Wang-Landau algorithm

In multicanonical sampling, we are given a state space \mathcal{X} and a probability measure π . \mathcal{X} is then partitioned as $\mathcal{X} = \cup \mathcal{X}_i$, where $\mathcal{X}_i \cap \mathcal{X}_j = \phi$ and π is re-weighted in each component \mathcal{X}_i . An abstract way to do the same and much more is the following. We start with $(\mathcal{X}_i, \mathcal{B}_i, \lambda_i)$ $i = 1, \dots, d$, a finite family of measure spaces where λ_i is a σ -finite measure. We introduce the union space $\mathcal{X} = \cup_{i=1}^d \mathcal{X}_i \times \{i\}$. We equip \mathcal{X} with the σ -algebra \mathcal{B} generated by $\{(A_i, i), i \in \{1, \dots, d\}, A_i \in \mathcal{B}_i\}$ and the measure λ satisfying $\lambda(A, i) = \lambda_i(A) \mathbf{1}_{\mathcal{B}_i}(A)$. Let $h_i : \mathcal{X}_i \rightarrow \mathbb{R}$ be a non-negative measurable function and define $\theta^*(i) = \int_{\mathcal{X}_i} h_i(x) \lambda_i(dx) / Z$ where $Z = \sum_{i=1}^d \int_{\mathcal{X}_i} h_i(x) \lambda_i(dx)$. We assume that $\theta^*(i) > 0$ for all $i = 1, \dots, d$ and consider the following probability measure on $(\mathcal{X}, \mathcal{B})$:

$$\pi^*(dx, i) \propto \frac{h_i(x)}{\theta^*(i)} \mathbf{1}_{\mathcal{X}_i}(x) \lambda_i(dx). \quad (1)$$

Our objective is to sample from π^* . The problem of sampling from such a distribution arises in a number of different Monte Carlo strategies. For example, and as explained above, if π is a probability measure of interest on some space $(\mathcal{X}, \mathcal{B}, \lambda)$, we can partition \mathcal{X} along the energy function $-\log(\pi)$ and re-weight π by $\pi(\mathcal{X}_i)$ in each component \mathcal{X}_i . The sampling problem then becomes of the form (1). This powerful strategy appeared first in the Physics literature as *multicanonical sampling* ([9]). This is discussed in some more details in Section 3.1.

Sampling from (1) also arises naturally when optimizing the *simulated tempering* algorithm ([18], [12]). In simulated tempering, the states space \mathcal{X} is not partitioned but, instead, some auxiliary distributions π_2, \dots, π_d are introduced (take $\pi_1 = \pi$). These distributions are chosen close to π but easier to sample from. For good performances, one typically imposes that all the distributions have the same weight. Taking each probability space $(\mathcal{X}, \mathcal{B}, \pi_i)$ as a component in the formalism above leads to a sampling problem of the form (1). Multicanonical sampling and simulated tempering have been combined in [5] giving an algorithm which can also be framed as (1). Sampling from (1) can also be an efficient strategy to improve on trans-dimensional MCMC samplers for Bayesian inference with model uncertainty. This is detailed in Section 3.3.

The main obstacle in sampling from π^* is that the normalizing constants θ^* are not known. The contribution of the Wang-Landau algorithm ([21]) is an efficient algorithm that simultaneously estimates θ^* and sample from π^* . The algorithm was introduced in a discrete setting with the π^* being uniform in i . In this work we extend the algorithm to general state spaces and to arbitrary probability measures. To carry on the discussion in our general framework, we introduce the family of probability measures $\{\pi_\theta, \theta \in (0, \infty)^d\}$ on $(\mathcal{X}, \mathcal{B}, \lambda)$ defined by:

$$\pi_\theta(dx, i) \propto \frac{h_i(x)}{\theta(i)} \mathbf{1}_{\mathcal{X}_i(x)} \lambda_i(dx). \quad (2)$$

We assume that for all $\theta \in (0, \infty)^d$, we have at our disposal a transition kernel P_θ on $(\mathcal{X}, \mathcal{B})$ with invariant distribution π_θ . Note that π_θ and P_θ remain unchanged if we multiply the vector θ by a positive constant. How to build such Markov chain P_θ typically depends on the particular instance of the algorithm. We give some examples later.

The structure of the WL algorithm is as follows. We start out with some initial value $(X_0, I_0) \in \mathcal{X}$, $\phi_0 \in (0, \infty)^d$ and set $\theta_0(i) = \phi_0(i) / \sum_{j=1}^d \phi_0(j)$, $i = 1, \dots, d$. Here θ_0 serves as an initial guess of θ^* . At iteration $n+1$, we generate (X_{n+1}, I_{n+1}) by sampling from $P_{\phi_n}(X_n, I_n; \cdot)$ and update ϕ_n to ϕ_{n+1} , which is used to form $\theta_{n+1}(i) = \phi_{n+1}(i) / \sum_j \phi_{n+1}(j)$. The updating rule for ϕ_n is fairly simple. For $i \in \{1, \dots, d\}$, if $X_{n+1} \in \mathcal{X}_i$ (equivalently, if $I_{n+1} = i$), then $\phi_{n+1}(i) = \phi_n(i)(1 + \rho)$ for some $\rho > 0$; otherwise $\phi_{n+1}(i) = \phi_n(i)$. This leads to a first version of the WL algorithm.

Algorithm 2.1 (The Wang-Landau algorithm I). *Let $\{\rho_n\}$ be a sequence of decreasing positive numbers. Let $(X_0, I_0) \in \mathcal{X}$ be given. Let $\phi_0 \in \mathbb{R}^d$ be such that $\phi_0(i) > 0$ and set $\theta_0(i) = \phi_0(i) / \sum_j \phi_0(j)$, $i = 1, \dots, d$. At some time $n \geq 0$, given $(X_n, I_n) \in \mathcal{X}$, $\phi_n \in \mathbb{R}^d$, $\theta_n \in \mathbb{R}^d$:*

- (i) *Sample $(X_{n+1}, I_{n+1}) \sim P_{\theta_n}(X_n, I_n; \cdot)$.*
- (ii) *For $i = 1, \dots, d$, set $\phi_{n+1}(i) = \phi_n(i) \left(1 + \rho_n \mathbf{1}_{\{I_{n+1}=i\}}\right)$ and $\theta_{n+1}(i) = \phi_{n+1}(i) / \sum_j \phi_{n+1}(j)$.*

It remains to choose the sequence $\{\rho_n\}$. As we show below, $\{\theta_n\}$ as defined by Algorithm 2.1 is a stochastic approximation process driven by $\{(X_n, I_n)\}$. The general guidelines in the literature to choose $\{\rho_n\}$ are: $\rho_n > 0$, $\sum \rho_n = \infty$ and $\sum \rho_n^{1+\varepsilon} < \infty$ for some $\varepsilon > 0$, often $\varepsilon = 1$. The typical choice is $\rho_n \propto n^{-1}$. In practice, more careful choices are often necessary for good performances. To the best of knowledge, there is no general, satisfactory way of choosing the step-size in stochastic approximation. Interestingly, Wang-Landau came up with a clever, adaptive way of choosing $\{\rho_n\}$

which works very well in practice. We describe their approach next, again in more probabilistic terms.

Let $v_{n,k}(i)$ denote the proportion of visits to $\mathcal{X}_i \times \{i\}$ between times $n + 1$ and k . That is, $v_{k,n}(i) = 0$ for $k \leq n$ and for $k \geq n + 1$, $v_{n,k}(i) = \frac{1}{k-n} \sum_{j=n+1}^k \mathbf{1}\{I_j = i\}$. Let $c \in (0, 1)$ be a parameter to be specified by the user. We introduce two additional random sequences $\{\kappa_n\}$ and $\{a_n\}$. Initially, $\kappa_0 = 0$. For $n \geq 1$, define

$$\kappa_n = \inf \left\{ k > \kappa_{n-1} : \max_{1 \leq i \leq d} \left| v_{\kappa_{n-1}, k}(i) - \frac{1}{d} \right| \leq \frac{c}{d} \right\}, \quad (3)$$

with the usual convention that $\inf \emptyset = \infty$. We need another sequence $\{\gamma_n\}$ of positive decreasing numbers, representing “stepsizes”. Then, $\{a_n\}$ represents the index of the element of the sequence $\{\gamma_n\}$ used at time n : $a_0 = 0$, if $k = \kappa_j$ for some $j \geq 1$, then $a_k = a_{k-1} + 1$ otherwise $a_k = a_{k-1}$. In other words, we start Algorithm 2.1 with a step-size equal to γ_0 and keep using it until time κ_1 when all the components are visited equally well. Only then we change the step-size to γ_1 and keep it constant until time κ_2 etc... Combining this with Algorithm 2.1, we get the following.

Algorithm 2.2 (The Wang-Landau algorithm II). *Let $\{\gamma_n\}$ be a sequence of decreasing positive numbers. Let $(X_0, I_0) \in \mathcal{X}$ be given. Set $a_0 = 0$, $\kappa = 0$, $c \in (0, 1)$, $\phi_0 \in \mathbb{R}^d$ such that $\phi_0(i) > 0$ and $\theta_0(i) = \phi_0(i) / \sum_j \phi_0(j)$, $i = 1, \dots, d$. At some time $n \geq 0$, given $(X_n, I_n) \in \mathcal{X}$, $\phi_n \in \mathbb{R}^d$, $\theta_n \in \mathbb{R}^d$, a_n and κ :*

- (i) *Sample $(X_{n+1}, I_{n+1}) \sim P_{\theta_n}(X_n, I_n; \cdot)$.*
- (ii) *For $i = 1, \dots, d$, set $\phi_{n+1}(i) = \phi_n(i) \left(1 + \gamma_{a_n} \mathbf{1}_{\{I_{n+1}=i\}} \right)$ and $\theta_{n+1}(i) = \phi_{n+1}(i) / \sum_j \phi_{n+1}(j)$.*
- (iii) *If $\max_i \left| v_{\kappa, n+1}(i) - \frac{1}{d} \right| \leq c/d$ then set $\kappa = n + 1$ and $a_{n+1} = a_n + 1$, otherwise $a_{n+1} = a_n$.*

Remark 2.1. 1. The performances of Algorithm 2.2 depends very much on the choice of $\{\gamma_n\}$ and c . Theoretically we show that the choice $\gamma_n \propto n^{-1}$ guaranties the convergence of the algorithm. But in practice, this type of step-size can be overly slow. The user might then consider instead $\gamma_n \propto a^{-n}$ ($a > 1$) originally proposed by Wang-Landau. But we were not able to obtain the convergence of the algorithm for summable step-sizes. A good compromise is to start the sampler with $\gamma_n \propto a^{-n}$ until $\gamma_n < \varepsilon$ (e.g. $\varepsilon = 10^{-5}$) and then switch to $\gamma_n = n^{-1}$. There is a bias-variance trade-off involved in the choice of c . For c close

to 0, Algorithm 2.2 will have a low bias (in estimating θ^*) but a high variance. Such values of c are more suitable for $\gamma_n \propto a^{-n}$. Whereas for larger values of c , the bias will be high with a low variance. For $c = d - 1$, we get a standard stochastic approximation algorithm with deterministic step-size for which a step-size that sums to ∞ is necessary for convergence. We found empirically from our simulations that c in the range $0.4 - 0.2$ yields reasonably good samplers.

2. In the actual implementation of the algorithm, it is not necessary to re-normalize ϕ_n into θ_n as in (ii). In fact, for computational stability, we recommend carrying out the recursion on a logarithmic scale: $\log \phi_n(i) = \log \phi_{n-1}(i) + \log(1 + \gamma_{a_{n-1}})\mathbf{1}_{\{I_n=i\}}$.
3. Another interesting feature of Algorithm 2.2 is that Step (iii) can serve as a stopping rule: we stop the simulation when γ_{a_n} get smaller than some pre-specified value.
4. Under some regularity conditions, if $f : \mathcal{X} \rightarrow \mathbb{R}$ is some function of interest that is π -integrable and $\{(X_n, I_n, \theta_n)\}$ is as described in Algorithm 2.2, we will show below that,

$$\frac{1}{n} \sum_{k=1}^n f(X_k, I_k) \rightarrow \pi^*(f), \text{ a.s. as } n \rightarrow \infty.$$

If we denote π_i the distribution on \mathcal{X}_i with density with respect to λ_i proportional to $h_i(x)\mathbf{1}_{\mathcal{X}_i}(x)$, we can estimate integrals with respect to π_i as well:

$$\frac{\sum_{k=1}^n f(X_k, I_k)\mathbf{1}_{\mathcal{X}_i}(X_k)}{\sum_{k=1}^n \mathbf{1}_{\mathcal{X}_i}(X_k)} \rightarrow \pi_i(f(\cdot, i)), \text{ a.s. as } n \rightarrow \infty.$$

Now if we denote π the distribution on \mathcal{X} whose density with respect to λ is proportional to $h_i(x)$ on \mathcal{X}_i , the ratio $\pi(x, i)/\pi^*(x, i)$ is $d\theta^*(i)$ and integrals with respect to π can also be computed by importance sampling:

$$\frac{d}{n} \sum_{k=1}^n f(X_k, I_k)\theta_k(I_k) \rightarrow \pi(f), \text{ a.s. as } n \rightarrow \infty.$$

Various methods for recycling Monte Carlo samples can be implemented as well. All these results follow from the strong law of large numbers of Theorem 4.1.

3. Some Applications

In this section, we detail briefly some applications of the general algorithm to multicanonical sampling and simulated tempering and trans-dimensional MCMC.

3.1. Multicanonical sampling

Multicanonical sampling is a powerful algorithm proposed by [9]. It holds the potential of improving on mixing times of classical MCMC algorithms. It fits naturally in the framework above. But the implementation can be tedious. Assume that we want to sample from a probability measure $\pi(dx) \propto h(x)\lambda(dx)$ on some probability space $(\Sigma, \mathcal{A}, \lambda)$. We use the energy function $E(x) = -\log(h(x))$ to build a d -component partition $(\mathcal{X}_i)_i$ of Σ . $\mathcal{X}_i = \{x \in \Sigma : E_{i-1} < E(x) \leq E_i\}$, where $-\infty \leq E_0 < E_1 < \dots < E_d \leq \infty$ are predefined values. Denote $\theta^*(i) = \pi(\mathcal{X}_i)$ and assume $\theta^*(i) > 0$. As above, we introduce the union space $\mathcal{X} = \bigcup \mathcal{X}_i \times \{i\}$. The idea of multicanonical sampling is to sample from π^* given by:

$$\pi^*(dx, i) \propto \frac{h(x)}{\theta^*(i)} \mathbf{1}_{\mathcal{X}_i}(x) \lambda(dx),$$

which is of the form (1). There is a simpler formulation of the algorithm. Since the component of the partition to which a point x belongs can be obtained from x itself, multicanonical sampling is equivalent to sampling from π^* on (Σ, \mathcal{A}) given by:

$$\pi^*(dx) \propto \sum_{i=1}^d \frac{h(x)}{\theta^*(i)} \mathbf{1}_{\mathcal{X}_i}(x) \lambda(dx), \quad (4)$$

and the union space formalism is not needed. After sampling from π^* in (4), a straightforward importance sampling estimate allows to recover π . The algorithm tries to break the barriers in the energy landscape of the distribution by re-weighting each component \mathcal{X}_i . Clearly, the success depends heavily on a good choice of the energy rings E_0, \dots, E_d . This typically requires some prior information on π or some pilot simulations. We point out that although the energy function E is a natural candidate to utilize to partition the space, the idea can be extended to other functions.

In the description of multicanonical sampling given above, taking Σ as a discrete space, π the uniform distribution on \mathcal{X} and $\mathcal{X}_e = \{x \in \Sigma : E(x) = e\}$, $e \in \{e \in \mathbb{R} : E(x) = e \text{ for some } x \in \Sigma\}$ yields the Wang-Landau algorithm of ([21]).

3.2. Simulated Tempering

The method can be applied to the simulated tempering of [18] and [12] by taking $(\mathcal{X}_i, \mathcal{B}_i, \lambda_i) \equiv (\mathcal{X}_1, \mathcal{B}_1, \lambda_1)$ and $h_i = h^{1/t_i}$, $1 = t_1 < t_2 < \dots < t_d$. Simulated tempering is a well-known Monte

Carlo strategy for sampling from difficult target distributions. Assume that the distribution of interest is $\pi_1(dx) \propto h(x)\lambda(dx)$. Typically for large temperature t , $h^{1/t}$ is a more well-behaved distribution for which faster mixing Markov chains can be built. In simulated tempering, we try to take advantage of these faster mixing chains, by targeting the distribution:

$$\pi_\theta(dx, i) = \frac{h^{1/t_i}(x)}{\theta(i)} \mathbf{1}_{\mathcal{X}_i}(x) \lambda_1(dx), \quad (5)$$

on the union space $(\mathcal{X}, \mathcal{B}, \lambda)$. A MCMC sampler P_θ with invariant distribution π_θ is readily designed. Typically, P_θ takes the form

$$P_\theta((x, i); A \times \{j\}) = B_{x,\theta}(i, j) P^{[j]}(x, A), \quad (6)$$

where $B_{x,\theta}$ is a transition kernel on $\{1, \dots, d\}$ with invariant distribution $(h_j(x)/\theta(j)) (\sum_i h_i(x)/\theta(i))^{-1}$ and $P^{[i]}$ a transition kernel on $(\mathcal{X}_i, \mathcal{B}_i)$ with invariant distribution proportional to $h_i(x)\lambda(dx)$. Typically, one takes $B_{x,\theta}(i, j) = (h_i(x)/\theta(i)) (\sum_i h_i(x)/\theta(i))^{-1}$. Another common choice is to take $B_{x,\theta}$ as a Metropolis-Kernel on $\{1, \dots, d\}$ with proposal $q(i, j)$. By standard importance sampling techniques, we can convert samples from the higher temperature distributions to also estimate π_1 . The method holds for any $\theta \in \mathbb{R}^d$, $\theta(i) > 0$. But the choice of θ can significantly impact the efficiency. Heuristically, it seems that, to improve on mixing, we need a θ that allows to sample from fast converging distributions (but close to π); but since π_1 is the distribution of interest, for statistical efficiency, we need a θ that favors π_1 . One easy way to resolve this trade-off is to choose θ such that all the distributions are equally visited. For this we need to choose $\theta(i) = \theta^*(i) \propto \int h_i(x)\lambda_1(dx)$ and sample from

$$\pi^*(dx, i) \propto \frac{h^{1/t_i}(x)}{\theta^*(i)} \mathbf{1}_{\mathcal{X}_i}(x) \lambda_1(dx),$$

which can be done with Algorithm 2.2.

Example 1. We compare a plain simulated tempering with weight $\theta(i) \equiv 1$ and the Wang-Landau simulated tempering described above for sampling from a multimodal bivariate Gaussian mixture distribution. The target distribution given below was taken from [17]

$$\pi(x) = \frac{1}{2\pi\sigma^2} \sum_{i=1}^{20} \omega_i \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu_i)' (x - \mu_i) \right\}, \quad (7)$$

where $\sigma = 0.1$ and $\omega_i \equiv 0.05$. The μ_i 's are listed in Table 1. The distribution is highly multimodal and it is clear that a plain Random Walk Metropolis algorithm for this distribution does not

mix in a reasonable time. Simulated tempering can be particularly efficient in such situations. We compare two strategies: a plain simulated tempering where $\theta(i) \equiv 1$ in (5) and the WL adaptation of simulated tempering as described above. We use the temperature scale $1 < 7.7 < 31.6 < 100$.

Table 2 presents the mean squared errors (MSEs) of the two methods in estimating the first two moments of the two components of π . We can see that the WL version is about three-four times more efficient than the plain version in terms of MSE. The estimates are based on 30 independent replications of the samplers. We run each sampler for 100,000 iterations. In applying Algorithm 2.2, we use $\gamma_n = 1/n$ and $c = 0.3$

i	μ_{i1}	μ_{i2}	i	μ_{i1}	μ_{i2}	i	μ_{i1}	μ_{i2}	i	μ_{i1}	μ_{i2}
1	2.18	5.76	6	3.25	3.47	11	5.41	2.65	16	4.93	1.50
2	8.67	9.59	7	1.70	0.50	12	2.70	7.88	17	1.83	0.09
3	4.24	8.48	8	4.59	5.60	13	4.98	3.70	18	2.26	0.31
4	8.41	1.68	9	6.91	5.81	14	1.14	2.39	19	5.54	6.86
5	3.93	8.82	10	6.87	5.40	15	8.33	9.50	20	1.69	8.11

TABLE 1

The 20 means of the two-dimensional Gaussian mixture.

	$\mathbb{E}(X_1)$	$\mathbb{E}(X_2)$	$\mathbb{E}(X_1^2)$	$\mathbb{E}(X_2^2)$
Plain ST	0.113	0.132	11.201	12.501
WL-ST	0.029	0.041	2.818	4.023
Ratio	3.89	3.25	3.97	3.11

TABLE 2

Mean Square Errors of the plain and the WL simulated tempering algorithms. Based on 30 independent replications with 100,000 iterations of each sampler.

3.3. Application to trans-dimensional MCMC

It is often the case in Statistics that many alternative models are considered for the same data. One is then interested in issues like model comparison, model selection and averaging. Let $f(\text{Data}|k, x_k)$ be the likelihood of model k with parameter x_k . Assume that we have a finite number d of models and that $x_k \in (\mathcal{X}_k, \mathcal{B}_k, \lambda_k)$. Let $\mathcal{X} = \bigcup_{i=1}^d \mathcal{X}_i \times \{i\}$ be the union space equipped as above with the σ -algebra \mathcal{B} and the σ -finite measure λ . $(\mathcal{X}, \mathcal{B}, \lambda)$ is the natural space to consider when dealing both with model uncertainty and parameter estimation. In the Bayesian framework,

a prior density (with respect to λ) $p(x_k, k)$ in $(\mathcal{X}, \mathcal{B})$ is specified for (x_k, k) . The posterior distribution of (x_k, k) is therefore $\pi(x_k, k) \propto h_k(x_k) = f(\text{Data}|k, x_k)p(x_k, k)$. In this framework, one is often interested in the Bayes factor of model i to model j defined as $B_{ij} := \theta^*(i)p(j)/(\theta^*(j)p(i))$, where $\theta^*(i) \propto \int_{\mathcal{X}_i} \pi(x_i, i)\lambda_i(dx_i)$ and $p(i) = \int p(x_i, i)\lambda_i(dx_i)$. Trans-dimensional MCMC is a set of specialized MCMC algorithms to sample from distributions like π defined on spaces of variable dimensions. The reversible-jump algorithm of Green ([14]) is the most popular such sampler.

In the spirit of the WL algorithm, an alternative to sampling directly from π is to sample from the distribution

$$\pi^*(dx_i, i) \propto \frac{h_i(x_i)}{\theta^*(i)} \mathbf{1}_{\mathcal{X}_i}(x) \lambda_i(dx_i). \quad (8)$$

By such re-weighting, we give the same posterior weight to all the models. The WL algorithm then offers an effective strategy to sample from π^* and we recover π by importance sampling. This strategy can improve on the mixing of the sampler

Example 2. We set $\mathcal{X}_i = \mathbb{R}^i$ for $i = 1, \dots, 20$ and consider the following rather trivial trans-dimensional target distribution:

$$\pi(x_i, i) \propto a_i^{-1} e^{-\frac{1}{2}|x_i|^2},$$

where we let $a_i = 1$ for $i \neq 4$, and $a_4 = (2\pi)^{-16/2}$. In this distribution, $x_i \in \mathbb{R}^i$, the i -dimensional Euclidean space and $\pi(x_i, i)$ restricted to \mathbb{R}^i is proportional to the standard normal distribution. We are interested in the marginal distribution $p(i)$ of i . This distribution, as shown in Figure 1, is bimodal with modes at 5 and 20.

We pretend that this distribution is intractable and sample from it using a Birth-and-Death Reversible-Jump MCMC. For the fixed-dimensional move, we use a Random Walk Metropolis kernel with a Gaussian proposal with covariance matrix $\sigma_p I_i$, $\sigma_p = 0.1$. We implement a Birth-and-Death move for the trans-dimensional jump. Given (x, i) , we randomly select $j \in \{i-1, i+1\}$ with respective probability $\omega_{i,i-1}, \omega_{i,i+1}$. We choose $\omega_{i,i+1} = 1/2$ with the usual correction at the boundaries. If $j = i+1$, we proposal $y = (x_i, u)$ where $u \sim N(0, \sigma^2)$ with $\sigma = 0.1$. We accept (y, j) with probability $\min(1, A)$ where

$$A = \frac{\pi(y, j) \omega_{ji}}{\pi(x, i) \omega_{ij}} \sqrt{2\pi\sigma} e^{-\frac{1}{2\sigma^2}u^2}.$$

Similarly if $j = i-1$, we write $x = (y, u')$ with $u' \in \mathbb{R}$ and propose (y, j) . This value is then

accepted with probability $\min(1, A)$ with

$$A = \frac{\pi(y, j) \omega_{ji}}{\pi(x, i) \omega_{ij}} \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2\sigma^2} w^2}.$$

This vanilla RJMCMC sampler fails to sample from π . Depending on its starting point, the sampler typically found one of the two modes and got stuck around that mode even after 10 millions iterations (see Fig 2a). In contrast, the WL algorithm provided a reasonable estimate of the distribution in only 2 millions iterations. In this example, computationally each WL iteration step costs roughly about that for 1.2 step of the vanilla sampler. For the WL approach we use $c = 0.4$ and $\gamma_n = 2^{-n}$ until 10^{-4} before switching to $\gamma_n = n^{-1}$.

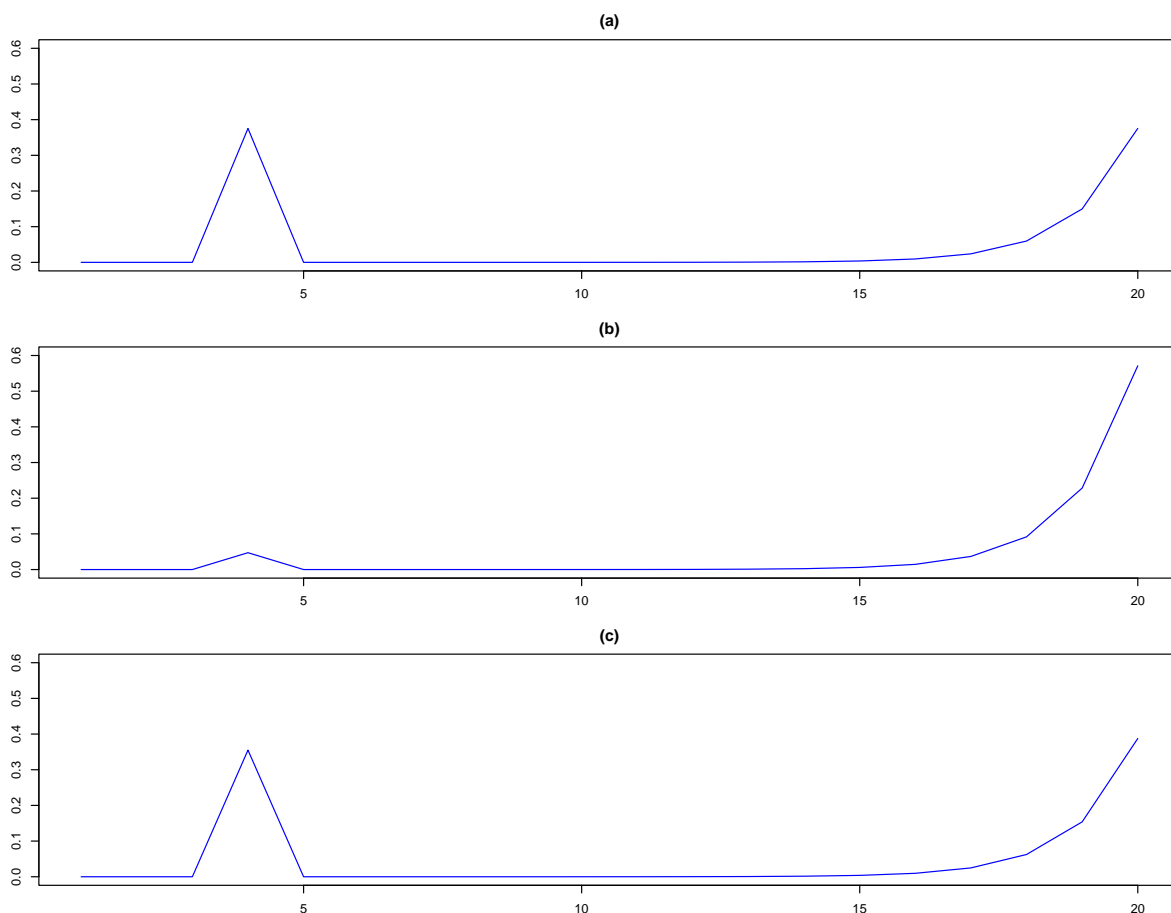


Figure 1: Marginal posterior distribution of models. (a) estimate from the plain RJMCMC; (b) estimate from the WL-RJMCMC ; (c) True posterior distribution. Estimates are based on 10×10^6 iterations for the plain RJMCMC and 2×10^6 iterations for the WL-RJMCMC.

4. Some theoretical results

We look at some theoretical aspects of the algorithm. We investigate the convergence of θ_n to θ^* and a strong law of large numbers for $\{(X_n, I_n)\}$. The difficulty is in proving that the algorithm is *stable* in the following sense.

Definition 4.1. Let $v_n(i)$ be the occupation measure of \mathcal{X}_i by time n : $v_n(i) = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{\{I_k=i\}}$. The Wang-Landau algorithm is said to be stable if

$$\max_{i,j} \limsup_{n \rightarrow \infty} n(v_n(i) - v_n(j)) < \infty, \text{ a.s.} \quad (9)$$

This is an essential property of the algorithm. When the algorithm is *stable*, we will show that all the stopping times κ_l defined in Algorithm 2.2 are finite and the step-size γ_{a_n} will gradually converges to 0. Moreover (Theorem 4.1 below) on a path where the algorithm is stable, $\theta_n \rightarrow \theta^*$ and a strong law of large numbers hold for $\frac{1}{n} \sum_{k=1}^n f(X_k, I_k)$. Then, we will derive some verifiable conditions under which the algorithm is shown to be stable. To maintain the flow of ideas, some of the proofs are postponed to Section 6.

4.1. Ergodicity

Let $\Theta = \{\theta \in \mathbb{R}^d : \sum_{i=1}^d \theta(i) = 1, \theta(i) \in (0, 1), i = 1, \dots, d\}$ and $((X_0, I_0), \theta_0) \in \mathcal{X} \times \Theta$ be the initial state of the algorithm. This initial state will be considered fixed but arbitrary. Let $\{\gamma_n\}$ be the step-size sequence. Let \Pr be the distribution of the process $\{X_n, I_n, \theta_n\}$ started at (X_0, I_0, θ_0) with step-size sequence $\{\gamma_n\}$ and denote \mathbb{E} the expectation with respect to \Pr . To simplify the notations, we will omit to make explicit the dependence of \Pr on (X_0, I_0, θ_0) and $\{\gamma_n\}$. All statements made almost surely will be with respect to \Pr . For $\theta \in \mathbb{R}^d$, $|\theta|$ denotes the Euclidean norm of θ . For any $\varepsilon \in (0, \varepsilon_*)$, where $\varepsilon_* = \min_i \theta^*(i)$, define $\Theta_\varepsilon = \{\theta \in \Theta, \theta(i) \geq \varepsilon, i = 1, \dots, d\}$. Our main assumption asserts that the family $\{P_\theta, \theta \in \Theta_\varepsilon\}$ is Lipschitz and uniformly V -ergodic. See e.g. [2] for some examples of MCMC samplers where these assumptions hold. Before stating these assumptions, we need some notations. For any function $f : \mathcal{X} \rightarrow \mathbb{R}$ and $W : \mathcal{X} \rightarrow [1, \infty)$ we denote $|f|_W := \sup_{x \in \mathcal{X}} \frac{|f(x)|}{W(x)}$, and introduce the set of W -bounded functions $L_W := \{f \text{ meas.}, f : \mathcal{X} \rightarrow \mathbb{R}, |f|_W < \infty\}$. A transition kernel on $(\mathcal{X}, \mathcal{B})$ operates on measurable real-valued functions f as $Pf(x) = \int P(x, dy)f(y)$ and the product of two transition kernels P_1

and P_2 is the transition kernel defined as $P_1 P_2(x, A) = \int P_1(x, dy) P_2(y, A)$. For two transition kernels P_1 and P_2 we define $\|P_1 - P_2\|_W$, the W -distance between P_1 and P_2 as

$$\|P_1 - P_2\|_W := \sup_{|f|_W \leq 1} |P_1 f - P_2 f|_W .$$

(A1) We assume the existence of a measurable function $V : \mathcal{X} \rightarrow [1, \infty)$, a set $C \subset \mathcal{X}$, a probability measure ν on $(\mathcal{X}, \mathcal{B})$ such that $\nu(C) > 0$ with the following property. For all $\varepsilon \in (0, \varepsilon_*)$ we can find constants $\lambda_\varepsilon \in (0, 1)$, $b_\varepsilon \in [0, \infty)$, $\beta_\varepsilon \in (0, 1]$ and integer $n_{0, \varepsilon}$ such that:

$$\inf_{\theta \in \Theta_\varepsilon} P_\theta^{n_{0, \varepsilon}}(x, A) \geq \beta_\varepsilon \nu(A) \mathbf{1}_C(x), \quad x \in \mathcal{X}, A \in \mathcal{B}; \quad (10)$$

and

$$\sup_{\theta \in \Theta_\varepsilon} P_\theta V(x) \leq \lambda_\varepsilon V(x) + b_\varepsilon \mathbf{1}_C(x), \quad x \in \mathcal{X}. \quad (11)$$

The inequality (11) of (A1) is the so-called drift condition and (10) is the so-called minorization condition.

Next, we assume that P_θ is Lipschitz as a function of θ .

(A2) For all $\alpha \in [0, 1]$, for all $\varepsilon \in (0, \varepsilon_*)$, there exists $K = K(\alpha, \varepsilon) < \infty$ such that for all $\theta, \theta' \in \Theta_\varepsilon$

$$\|P_\theta - P_{\theta'}\|_{V^\alpha} \leq K |\theta - \theta'|, \quad (12)$$

where V is defined in (A1).

Finally we assume that the step-size sequence is well-behaved:

(A3) $\{\gamma_n\}$ is non-increasing, $\gamma_n > 0$, $\sum \gamma_n = \infty$ and $\gamma_n = O(n^{-1})$ as $n \rightarrow \infty$.

For $B > 0$, we introduce the following stopping time:

$$\tau(B) := \inf\{k \geq 0 : \max_{i,j} k(v_k(i) - v_k(j)) > B\}, \quad (13)$$

with the usual convention that $\inf \emptyset = \infty$. $\tau(B)$ is the first time where one component will accumulate B visits or more than some other component. Thus Definition 4.1 is precisely equivalent to $\tau(B) = \infty$ for some B . With this in mind, the next results says essentially that if the algorithm is *stable* and (A1-3) hold then it is ergodic.

Theorem 4.1. *Assume (A1-3). Let $B > 0$ be given. Then:*

(i)

$$|\theta_n - \theta^*| \mathbf{1}_{\{\tau(B) > n\}} \rightarrow 0, \text{ a.s. as } n \rightarrow \infty. \quad (14)$$

(ii) For any function $f \in L_{V^{1/2}}$, denoting $\bar{f} = f - \pi^*(f)$, we have:

$$\frac{1}{n} \sum_{k=1}^n \bar{f}(X_k, I_k) \mathbf{1}_{\{\tau(B) > k-1\}} \rightarrow 0, \text{ a.s. as } n \rightarrow \infty. \quad (15)$$

Proof. See Section 6.2 □

4.2. Checking (A1-2)

We impose the drift condition and the minorization condition (A1) uniformly for $\theta \in \Theta_\varepsilon$, not uniformly for $\theta \in \Theta$. This is an important point. Indeed and as we will see now, a drift and minorization condition uniformly-in- θ for $\theta \in \Theta_\varepsilon$ is almost always true as soon as one P_θ satisfies these conditions (Proposition 4.1). Whereas a minorization and drift condition uniformly-in- θ for $\theta \in \Theta$ is almost never true.

Indeed, suppose that each P_θ is a Metropolis-Hastings kernel with invariant distribution π_θ and proposal kernel Q . That is

$$P_\theta f(x, i) = M_\theta f(x, i) + r_\theta(x, i) f(x, i),$$

where

$$M_\theta f(x, i) = \sum_{j=1}^d \int_{\mathcal{X}_j} \min \left(1, \frac{\theta(i)}{\theta(j)} R(y, j; x, i) \right) f(y, j) Q(x, i; dy, j)$$

and

$$r_\theta(x, i) = 1 - \sum_{j=1}^d \int_{\mathcal{X}_j} \min \left(1, \frac{\theta(i)}{\theta(j)} R(y, j; x, i) \right) Q(x, i; dy, j),$$

and $R(x, i; y, j)$ the Radon-Nikodym density of $\pi(dx, i)Q(x, i; dy, j)$ with respect to $\pi(dy, j)Q(y, j; dx, i)$.

With $\theta = (1, \dots, 1)$, we use π (resp. P and M) to denote π_θ (resp. P_θ and M_θ). The next result states that we only need to check that P is geometrically ergodic to obtain (A1-2).

Proposition 4.1. *Suppose that existence of a measurable function $V : \mathcal{X} \rightarrow [1, \infty)$, a set $C \subset \mathcal{X}$, a probability measure ν on $(\mathcal{X}, \mathcal{B})$ such that $\nu(C) > 0$; constants $\lambda \in (0, 1)$, $b \in [0, \infty)$, $\beta \in (0, 1]$ and finite integer n_0 such that:*

$$M^{n_0}(x, A) \geq \beta \nu(A) \mathbf{1}_C(x), \quad x \in \mathcal{X}, A \in \mathcal{B};$$

and

$$PV(x) \leq \lambda V(x) + b \mathbf{1}_C(x), \quad x \in \mathcal{X}.$$

Then (A1-2) hold.

Proof. See Section 6.1. □

4.3. Stability

Theorem 4.1 asserts that under (A1-3), the WL algorithm will converge to the right limit on *stable* paths. This opens the question of checking the stability of the algorithm. The stability condition is difficult to check in general. The next theorem gives some easily checked conditions under which the algorithm is stable.

Theorem 4.2. *The WL algorithm is stable under either of the following two conditions.*

- (a) *There exist $\varepsilon \in (0, 1)$, $K \in (0, \infty)$ and integer $n_0 \geq 0$ such that for any $i, j \in \{1, \dots, d\}$ and $\theta \in \mathbb{R}^d$, $\theta(i)/\theta(j) > K$ implies that $P_\theta^{n_0}((x, i), \mathcal{X}_j \times \{j\}) \geq \varepsilon$ for all $x \in \mathcal{X}_i$.*
- (b) *There exists $\varepsilon \in (0, 1)$, such that for any $j \in \{1, \dots, d\}$ and $\theta \in \mathbb{R}^d$, $\theta(j) \leq \min_{1 \leq i \leq d} \theta(i)$ implies $P_\theta((x, i), \mathcal{X}_j \times \{j\}) \geq \varepsilon$ for all $x \notin \mathcal{X}_i$.*

Proof. See Section 6.3. □

4.4. Application to Multicanonical sampling

Consider the multicanonical sampling of Section 3.1. Suppose that P_θ is the *Independence-Sampler* with proposal distribution $Q(dx) = q(x)\lambda(dx)$ and invariant distribution $\pi_\theta(dx) \propto \sum_{i=1}^d (h(x)/\theta(i)) \mathbf{1}_{\mathcal{X}_i}(x) \lambda(dx)$. Assume that the function $\omega(x) \propto h(x)/q(x)$ is bounded with supremum ω_0 . Then clearly, for all $x \in \mathcal{X}_i$,

$$\begin{aligned} P_\theta(x, \mathcal{X}_j) &\geq \min\left(1, \frac{\theta(i)}{\theta(j)}\right) \int_{\mathcal{X}_j} \min\left(1, \frac{\omega(y)}{\omega_0}\right) Q(dy) \\ &\geq \varepsilon_j, \end{aligned}$$

as soon as $\theta(i) \geq \theta(j)$, taking $\varepsilon_j = \int_{\mathcal{X}_j} \min\left(1, \frac{\omega(y)}{\omega_0}\right) Q(dy) > 0$ (since $\pi(\mathcal{X}_i) > 0$). Thus for any $i, j \in \{1, \dots, d\}$, $i \rightsquigarrow j$ and by Theorem 4.2, the WL algorithm is *stable*. Now, since ω is

bounded, each P_θ satisfies a drift condition and a minorization condition which implies (A1-2) by Proposition 4.1.

Similarly, if P_θ is a Random Walk Metropolis with proposal kernel $q(y-x)$ we have:

$$P_\theta(x, \mathcal{X}_j) \geq \min\left(1, \frac{\theta(i)}{\theta(j)}\right) \int_{\mathcal{X}_j} \min\left(1, \frac{\pi(y)}{\pi(x)}\right) q(y-x) dy.$$

It follows that if \mathcal{X} is compact and π, q positive and continuous, then $i \rightsquigarrow j$ for all i, j . Under the same assumption, (A1-2) also hold. We can conclude that:

Corollary 4.1. *In the case of multicanonical sampling of Section 3.1. Assume either (i) or (ii):*

- (i) P_θ is an independent-Metropolis sampler with proposal distribution $Q(dx) = q(x)\lambda(dx)$ and $\omega \propto h/q$ is bounded.
- (ii) P_θ is a RWM sampler with proposal $q(y-x)$; \mathcal{X} is compact and π and q are positive and continuous.

Then the algorithm is stable and under (A3), we have:

$$|\theta_n - \theta^*| \rightarrow 0; \quad \text{and} \quad \frac{1}{n} \sum_{k=1}^n f(X_k) \rightarrow \pi^*(f) \quad \text{a.s. as } n \rightarrow \infty,$$

for any bounded measurable function f .

4.5. Application to Simulated tempering

Theorems 4.1 and 4.2 can also be applied to the WL version of simulated tempering as described in Section 3.2. We consider the simulated tempering algorithm with kernel $P_\theta((x, i); A \times \{j\}) = B_{x,\theta}(i, j)P^{[j]}(x, A)$ where $B_{x,\theta}(i, j) \propto h_j(x)/\theta(j)$ and $P^{[j]}$ is a transition kernel (not necessarily Metropolis-Hastings) with invariant distribution h_j . The following corollary is easily proved and is left to the reader.

Corollary 4.2. *Suppose that \mathcal{X}_1 is compact, h positive and continuous and $\{\gamma_n\}$ satisfies (A3). Suppose that there exist $\varepsilon > 0$, a probability measure ν and an integer n_0 such that $\nu(\mathcal{X}_i) > 0$ and $(P^{[i]})^{n_0}(x, A) \geq \varepsilon\nu(A)$, $i \in \{1, \dots, d\}$. Then*

$$|\theta_n - \theta^*| \rightarrow 0; \quad \text{and} \quad \frac{1}{n} \sum_{k=1}^n f(X_k, I_k) \rightarrow \pi^*(f) \quad \text{a.s. as } n \rightarrow \infty,$$

for any bounded measurable function f .

5. Discussion and open problems

In this paper, we propose an extension of the WL algorithm to general states spaces. The WL algorithm differs from other adaptive Markov Chain Monte Carlo algorithms based on stochastic approximation by the adaptive nature of its step-size. We have shown through examples that the algorithm can be used effectively to improve on simulated tempering and trans-dimensional MCMC algorithms. We have also studied the asymptotic behavior of the WL algorithm. We have shown that on stable sample paths and with an appropriate step-size, θ_n converges to θ^* and a strong law of large numbers hold. Finally, when the state space is compact, we have shown that in most cases, the algorithm is stable and the aforementioned limit results apply.

Two main questions have remained unanswered. The first question concerns the stability of the algorithm in unbounded state spaces. Secondly, in order to exploit the full potential of the algorithm, a more precise understanding of its efficiency is needed. In particular we need to understand how the rate of convergence and the asymptotic variances of the algorithm are related to the parameter c and step-size γ_n . We are currently investigating some of these questions.

6. Proofs

The techniques used here can be found in various forms in [8], [11], [3]. Throughout, C_ε denotes a generic constant whose value can be different from one equation to another. The key result in the proof of Theorem 4.1 is Lemma 6.6 which states that the weighted sum of the noise process in the stochastic approximation followed by $\{\theta_n\}$ is summable.

6.1. Proof of Proposition 4.1

Lemma 6.1. *Under the conditions of Proposition 4.1, (A2) hold.*

Proof. Let $\varepsilon \in (0, \varepsilon^*)$ and $\alpha \in (0, 1]$. For $\theta \in \Theta_\varepsilon$, $|f| \leq V^\alpha$ and $(x, i) \in \mathcal{X}$, we have:

$$P_\theta f(x, i) = M_\theta f(x, i) + f(x, i) (1 - M_\theta \mathbf{e}(x, i)),$$

where $\mathbf{e}(x, i) \equiv 1$, and

$$M_\theta f(x, i) = \sum_{j=1}^d \int_{\mathcal{X}_j} \min \left(1, \frac{\theta(i)}{\theta(j)} r(y, j; x, i) \right) f(y, j) Q(x, i; dy, j).$$

As a consequence, for $\theta_2, \theta_1 \in \Theta_\varepsilon$:

$$\|P_{\theta_2} - P_{\theta_1}\|_{V^\alpha} \leq \|M_{\theta_2} - M_{\theta_1}\|_{V^\alpha} + \|M_{\theta_2} - M_{\theta_1}\|_{TV},$$

where $\|P\|_{TV} = \|P\|_V$ with $V \equiv 1$. For $(x, i) \in \mathcal{X}$, the function $\theta \rightarrow M_\theta f(x, i)$ is differentiable and:

$$\sum_{j=1}^d \left| \frac{\partial}{\partial \theta_j} M_\theta f(x, i) \right| \leq C_\varepsilon \sum_{j=1}^d M |f| (x, i).$$

(A2) then follows by the mean-value theorem and the drift condition on P . \square

Next we show that for any $\varepsilon \in (0, \varepsilon^*)$, the family $(P_\theta)_{\theta \in \Theta_\varepsilon}$ satisfies a uniform (in θ) drift condition.

Lemma 6.2. *Under the conditions of Proposition 4.1, (A1) hold.*

Proof. The minorization is immediate. Since $\min(1, ab) \geq \min(1, a) \min(1, b)$ for $a, b > 0$, $P_\theta(x, i; A) \geq M_\theta(x, i; A) \geq (1/\varepsilon)M_0(x, i; A) \geq \beta/\varepsilon\nu(A)\mathbf{1}_C(x, i)$.

The argument to show the uniform drift is fairly simple and we only sketch it. Let $B_{\theta, r} = \Theta_\varepsilon \cap B(\theta, r)$ the open ball of Θ_ε with center $\theta \in \Theta_\varepsilon$ and radius $r > 0$. It follows from Lemma 6.1 that by choosing $r > 0$ small enough, if P_θ satisfies a drift condition toward a small set C , then the family $\{P_{\theta'}, \theta' \in B_{\theta, r}\}$ satisfies a uniform (in θ') drift condition toward C . Therefore, starting from P , we can find an open coverage of Θ_ε by the $B_{\theta, r}$ such that on each such $B_{\theta, r}$ a uniform drift toward C hold. Since Θ_ε is compact it admits a finite coverage by the $B_{\theta, r}$ and using the maximum of the constants of the drift condition toward C for each such ball, we get a uniform drift toward C for $\{P_\theta, \theta \in \Theta_\varepsilon\}$. \square

6.2. Proof of Theorem 4.1

We start with some preliminary remarks. For any $\varepsilon \in (0, \varepsilon_*)$ and $\alpha \in (0, 1]$, it is known from Markov chains theory that (A1) implies the existence of $C_{\varepsilon, \alpha} < \infty$, $\rho_{\varepsilon, \alpha} \in (0, 1)$ and b_ε as in (A1) such that:

$$\sup_{\theta \in \Theta_\varepsilon} \|P_\theta^n - \pi_\theta\|_{V^\alpha} \leq C_{\varepsilon, \alpha} \rho_{\varepsilon, \alpha}^n, \quad (16)$$

and

$$\sup_{\theta \in \Theta_\varepsilon} \pi_\theta(V) \leq b_\varepsilon. \quad (17)$$

For a proof, see e.g. [7] and the references therein. Define $\xi(\varepsilon) = \inf\{k \geq 0 : \theta_k \notin \Theta_\varepsilon\}$. An easy calculation using (A1) gives that

$$\begin{aligned} \mathbb{E} \left[V(X_n, I_n) \mathbf{1}_{\{\xi(\varepsilon) > n\}} \right] &= \mathbb{E} [V(X_n, I_n) \mathbf{1}_{\Theta_\varepsilon}(\theta_0) \cdots \mathbf{1}_{\Theta_\varepsilon}(\theta_n)] \\ &\leq \lambda_\varepsilon^n V(X_0, I_0) + b_\varepsilon / (1 - \lambda_\varepsilon) \end{aligned}$$

from which we deduce:

$$\sup_n \mathbb{E} \left(V(X_n, I_n) \mathbf{1}_{\{\xi(\varepsilon) > n\}} \right) < \infty. \quad (18)$$

The proof of the theorem will be based on some nice properties of solutions of so-called Poisson equation. These solutions will allow us to obtain a martingale approximation to the process $\sum_{k=1}^n f(X_k, I_k)$. For $f \in L_{V^\alpha}$, $\alpha \in (0, 1]$ and $\theta \in \Theta_\varepsilon$, define the function

$$h_\theta = \sum_{k=0}^{\infty} P_\theta^k (f - \pi_\theta(f)). \quad (19)$$

h_θ solves the Poisson equation $f - \pi_\theta(f) = h_\theta - P_\theta h_\theta$. By (A1), h_θ exists and $h_\theta \in L_{V^\alpha}$. For $f \in L_{V^\alpha}$ and $\theta, \theta' \in \Theta_\varepsilon$, denoting $\bar{f}_\theta = f - \pi_\theta(f)$, we have

$$\begin{aligned} |\pi_\theta(f) - \pi_{\theta'}(f)| &= \left| \pi_\theta \left[P_\theta^k \bar{f}_{\theta'} \right] \right| \\ &= \left| \pi_\theta \left[P_{\theta'}^k (\bar{f}_{\theta'}) \right] + \sum_{j=1}^k \pi_\theta \left[P_\theta^{k-j} (P_\theta - P_{\theta'}) P_{\theta'}^{j-1} (\bar{f}_{\theta'}) \right] \right| \\ &\leq C_\varepsilon \left(\rho_\varepsilon^k + |\theta - \theta'| \sum_{j=1}^k \rho_\varepsilon^{j-1} \right), \end{aligned}$$

using (A1-2) from which we deduce that there exists a finite constant C_ε such that:

$$\sup_{f \in L_{V^\alpha}} |\pi_\theta(f) - \pi_{\theta'}(f)| \leq C_\varepsilon |\theta - \theta'|. \quad (20)$$

The constant C_ε is not necessarily the same from one equation to the other. Similarly, denoting \bar{P}_θ the operator $P_\theta - \pi_\theta$, we have

$$\begin{aligned} \left| P_\theta^k \bar{f}_\theta - P_{\theta'}^k \bar{f}_{\theta'} \right| &= \left| \bar{P}_\theta^k f - \bar{P}_{\theta'}^k f \right| \\ &= \left| \sum_{j=1}^k \bar{P}_\theta^{k-j} (P_\theta - P_{\theta'}) \bar{P}_{\theta'}^{j-1} f \right| \\ &\leq C_\varepsilon |\theta - \theta'| k \rho_\varepsilon^{k-1} V^\alpha, \end{aligned}$$

using (A1-2). This, together with (20) imply the existence of a finite constant C_ε such that for all $\alpha \in (0, 1]$, $\theta, \theta' \in \Theta_\varepsilon$:

$$|h_\theta - h_{\theta'}|_{V^\alpha} + |P_\theta h_\theta - P_{\theta'} h_{\theta'}|_{V^\alpha} \leq C_\varepsilon |\theta - \theta'|. \quad (21)$$

In our analysis, we mainly see $\{\theta_n\}$ as a stochastic approximation sequence. The recursion on $\{\theta_n\}$ writes:

$$\begin{aligned} \theta_{n+1}(i) &= \frac{\phi_{n+1}(i)}{\sum_{e=1}^d \phi_{n+1}(e)} \\ &= \frac{\phi_n(i) + \gamma_{a_n} \phi_n(i) \mathbf{1}_{\{I_{n+1}=i\}}}{\sum_{e=1}^d \phi_n(e) + \gamma_{a_n} \phi_n(I_{n+1})} \\ &= \theta_n(i) \frac{1 + \gamma_{a_n} \mathbf{1}_{\{I_{n+1}=i\}}}{1 + \gamma_{a_n} \theta_n(I_{n+1})} \\ &= \theta_n(i) + \gamma_{a_n} H_i(\theta_n, I_{n+1}) + \gamma_{a_n}^2 r_{i,n}(\theta_n, I_{n+1}), \end{aligned} \quad (22)$$

where $H_i(\theta, I) = \theta(i) (\mathbf{1}_{\{I=i\}} - \theta(I))$ and $r_{i,n}(\theta, I) = -\theta(i)\theta(I) \frac{\mathbf{1}_{\{I=i\}} - \theta(I)}{1 + \gamma_{a_n} \theta(I)}$.

The mean field function $h_i(\theta) = \pi_\theta(H_i)$ is

$$h_i(\theta) = \frac{\theta^*(i) - \theta(i)}{\sum_{j=1}^d \frac{\theta^*(j)}{\theta(j)}}. \quad (23)$$

Lemma 6.3. *Assume (A3). Let $B > 0$ be given and $c_0 = 2Bd/c$. Then on $\{\tau(B) > n\}$, $a_n \geq \lfloor n/c_0 \rfloor$. Moreover*

$$\sum \gamma_{a_n} = \infty, \quad \text{and} \quad \sum \gamma_{a_n}^2 \mathbf{1}_{\{\tau(B) > n\}} < \infty \quad \text{a.s.}$$

Proof. On $\{\tau(B) > n\}$, for any $k < k' \leq n$, and for any $i \in \{1, \dots, d\}$, $(v_{k,k'}(i) - 1/d) \leq 2B/(k' - k)$. Therefore if $lc_0 \leq n < \tau(B)$, $\kappa_l \leq n$. That is, $a_n \geq \lfloor n/c_0 \rfloor$ on $\{\tau(B) > n\}$.

Since $\{\gamma_n\}$ is non-increasing and $a_n \leq n$, $\sum \gamma_{a_n} \geq \sum \gamma_n = \infty$. On the other hand $\sum \gamma_{a_n}^2 \mathbf{1}_{\{\tau(B) > n\}} \leq Cc_0 \sum n^{-2} < \infty$. \square

For any $\varepsilon > 0$, we introduce the stopping time

$$\xi(\varepsilon) := \inf \{k \geq 0 : \theta_k \notin \Theta_\varepsilon\}.$$

We will need the following lemma whose proof is left to the reader.

Lemma 6.4. *Let $\{\gamma_n, n \geq 0\}$ be a non-increasing sequence of positive number and $\{v_n, n \geq 0\}$ a sequence of numbers such that $|\sum_{k=0}^N v_k| \leq B$ for all $N \geq 0$. Then $|\sum_{k=0}^N \gamma_k v_k| \leq \gamma_0 B$ for all $N \geq 0$.*

The following lemma relates $\tau(B)$ and $\xi(\varepsilon)$.

Lemma 6.5. *Assume (A3). For any $B > 0$, we can find $\varepsilon \in (0, \varepsilon^*)$ such that $\tau(B) \leq \xi(\varepsilon)$.*

Proof. Take $\varepsilon = \left(1 + (d-1)e^{B\gamma_0}\right)^{-1} > 0$. Without any loss, we assume that $\varepsilon < \varepsilon^*$ and $\varepsilon \leq \min_i \theta_0(i)$. We need to show that $\min_i \theta_n(i) > \varepsilon$ for all $n < \tau(B)$. But $\theta_n(i) = \left(1 + \sum_{j \neq i} \frac{\phi_n(j)}{\phi_n(i)}\right)^{-1}$. It is thus enough to show that $\phi_n(j)/\phi_n(i) \leq e^{B\gamma_0}$ for all $i \neq j$ and for any $n < \tau(B)$. But we have:

$$\frac{\phi_n(j)}{\phi_n(i)} = \exp\left(\sum_{p=0}^{\infty} \gamma_p (N_{\kappa_p, n \wedge \kappa_{p+1}}(j) - N_{\kappa_p, n \wedge \kappa_{p+1}}(i))\right),$$

where $N_{l,m}(i) = 0$ if $m \leq l$ and $N_{l,m}(i) = \sum_{q=l+1}^m \mathbf{1}_{\mathcal{X}_i}(X_q)$ otherwise ($N_{l,m}(i)$ is the number of visits to \mathcal{X}_i from time $l+1$ to m). For any $n < \tau(B)$, $|\sum_{p=0}^P (N_{\kappa_p, n \wedge \kappa_{p+1}}(i) - N_{\kappa_p, n \wedge \kappa_{p+1}}(j))| \leq B$ for all $P \geq 0$. Lemma 6.4 thus implies that $\phi_n(i)/\phi_n(j) \leq e^{\gamma_0 B}$. \square

Lemma 6.6. *Assume (A1-3). Let $B > 0$ be given. Let $\{\gamma'_n\}$ be a sequence that satisfies (A3) and such that $\sum \gamma_{\lfloor n/a \rfloor} \gamma'_{\lfloor n/a \rfloor} < \infty$ for all $a > 0$. For $\theta \in \Theta$, Let $H_\theta : \mathcal{X} \rightarrow \mathbb{R}$ be a measurable function such that $H_\theta \in L_{V^{1/2}}$. Then:*

$$\sum_{k=0}^{\infty} \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} [H_{\theta_k}(X_{k+1}, I_{k+1}) - \pi_{\theta_k}(H_{\theta_k})] < \infty, \quad a.s. \quad (24)$$

Proof. Let $\varepsilon > 0$ as in Lemma 6.5. From (A1), $h_\theta \in L_{V^{1/2}}$ exists that solves the Poisson equation $h_\theta - P_\theta h_\theta = H_\theta - \pi_\theta(H_\theta)$ for all $\theta \in \Theta_\varepsilon$. Using this we can write:

$$\sum_{k=0}^n \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} [H_{\theta_k}(X_{k+1}, I_{k+1}) - \pi_{\theta_k}(H_{\theta_k})] = \sum_{k=0}^n \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} \left(U_{k+1}^{(1)} + U_{k+1}^{(2)} + U_{k+1}^{(3)} \right),$$

where

$$\begin{aligned} U_{k+1}^{(1)} &= h_{\theta_k}(X_{k+1}, I_{k+1}) - P_{\theta_k} h_{\theta_k}(X_k, I_k), \\ U_{k+1}^{(2)} &= P_{\theta_k} h_{\theta_k}(X_k, I_k) - P_{\theta_{k+1}} h_{\theta_{k+1}}(X_{k+1}, I_{k+1}), \\ U_{k+1}^{(3)} &= P_{\theta_{k+1}} h_{\theta_{k+1}}(X_{k+1}, I_{k+1}) - P_{\theta_k} h_{\theta_k}(X_{k+1}, I_{k+1}). \end{aligned}$$

Clearly, $M_n = \sum_{k=0}^n \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} U_{k+1}^{(1)}$ is a martingale and using Lemma 6.3, (18) and since h_θ

and $P_\theta h_\theta \in L_{V^{1/2}}$, we have:

$$\begin{aligned} \mathbb{E} \left(M_n^2 \right) &\leq C_\varepsilon \sum_{k=0}^n \mathbb{E} \left[(\gamma'_{a_k})^2 \mathbf{1}_{\{\tau(B) > k\}} V(X_{k+1}, I_{k+1}) \right] \\ &\leq C_\varepsilon \sum_{k=0}^n (\gamma'_{[k/c_0]})^2 \mathbb{E} \left[\mathbf{1}_{\{\tau(B) > k\}} V(X_{k+1}, I_{k+1}) \right] \\ &\leq C_\varepsilon \sum_{k=0}^{\infty} (\gamma'_{[k/c_0]})^2 < \infty. \end{aligned}$$

By Doob's convergence theorem for martingales, $\sum_{k=0}^{\infty} \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} U_{k+1}^{(1)}$ is finite a.s.

On $\{\tau(B) = l\}$, $l < \infty$, $\sum_{k=0}^{\infty} \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} U_{k+1}^{(2)} = \sum_{k=0}^{l-1} \gamma'_{a_k} U_{k+1}^{(2)}$ which is finite almost surely.

On $\{\tau(B) = \infty\}$, we can write:

$$\begin{aligned} \mathbf{1}_{\{\tau(B) = \infty\}} \sum_{k=0}^n \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} U_{k+1}^{(2)} &= \gamma'_{a_0} P_{\theta_0} h_{\theta_0}(X_0, I_0) - \gamma'_{a_n} \mathbf{1}_{\{\tau(B) = \infty\}} P_{\theta_{n+1}} h_{\theta_{n+1}}(X_{n+1}, I_{n+1}) \\ &\quad + \sum_{k=0}^{n-1} (\gamma'_{a_{k+1}} - \gamma'_{a_k}) \mathbf{1}_{\{\tau(B) = \infty\}} P_{\theta_{k+1}} h_{\theta_{k+1}}(X_{k+1}, I_{k+1}). \end{aligned}$$

$\mathbb{E} \left[\left(\gamma'_{a_n} \mathbf{1}_{\{\tau(B) = \infty\}} P_{\theta_{n+1}} h_{\theta_{n+1}}(X_{n+1}, I_{n+1}) \right)^2 \right] \leq C_\varepsilon (\gamma'_{[n/c_0]})^2$ and $\sum (\gamma'_{[n/c_0]})^2 < \infty$ thus

$\gamma'_{a_n} \mathbf{1}_{\{\tau(B) = \infty\}} P_{\theta_{n+1}} h_{\theta_{n+1}}(X_{n+1}, I_{n+1})$ converges a.s. to 0.

$$\sum_{k=0}^{n-1} \left| (\gamma'_{a_{k+1}} - \gamma'_{a_k}) \mathbf{1}_{\{\tau(B) = \infty\}} P_{\theta_{k+1}} h_{\theta_{k+1}}(X_{k+1}, I_{k+1}) \right| \leq C_\varepsilon \sum_{k=0}^{\infty} (\gamma'_{a_k} - \gamma'_{a_{k+1}}) \mathbf{1}_{\{\tau(B) = \infty\}} V^{1/2}(X_{k+1}, I_{k+1}).$$

Since a_k only changes at the stopping times κ_i , we have

$$\begin{aligned} \mathbb{E} \left[\sum_{k=0}^{\infty} (\gamma'_{a_k} - \gamma'_{a_{k+1}}) \mathbf{1}_{\{\tau(B) = \infty\}} V^{1/2}(X_{k+1}, I_{k+1}) \right] &= \mathbb{E} \left[\sum_{k=1}^{\infty} (\gamma'_{k-1} - \gamma'_k) \mathbf{1}_{\{\tau(B) = \infty\}} V^{1/2}(X_{\kappa_k+1}, I_{\kappa_k+1}) \right] \\ &= \sum_{k=1}^{\infty} (\gamma'_{k-1} - \gamma'_k) \mathbb{E} \left[\mathbf{1}_{\{\tau(B) = \infty\}} V^{1/2}(X_{\kappa_k+1}, I_{\kappa_k+1}) \right] \\ &\leq C_\varepsilon \sum_{k=1}^{\infty} (\gamma'_{k-1} - \gamma'_k) \\ &\leq C_\varepsilon \gamma'_0. \end{aligned}$$

By Lebesgue's dominated convergence theorem, we can conclude that

$$\mathbb{E} \left[\left| \sum_{k=0}^{\infty} (\gamma'_{a_{k+1}} - \gamma'_{a_k}) \mathbf{1}_{\{\tau(B) = \infty\}} P_{\theta_{k+1}} h_{\theta_{k+1}}(X_{k+1}, I_{k+1}) \right| \right] < \infty,$$

This is sufficient to conclude that $\sum_{k=0}^{\infty} \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} U_{k+1}^{(2)}$ is finite almost surely.

Using (21):

$$\begin{aligned} \sum_{k=0}^n \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} \left| U_{k+1}^{(3)} \right| &\leq C_\varepsilon \sum_{k=0}^{\infty} \gamma'_{a_k} \gamma_{a_k} \mathbf{1}_{\{\tau(B) > k\}} V^{1/2}(X_{k+1}, I_{k+1}) \\ &\leq C_\varepsilon \sum_{k=0}^{\infty} \gamma'_{[n/c_0]} \gamma_{[n/c_0]} \mathbf{1}_{\{\tau(B) > k\}} V^{1/2}(X_{k+1}, I_{k+1}) \end{aligned}$$

and

$$E \left[\sum_{k=0}^{\infty} \gamma'_{[n/c_0]} \gamma_{[n/c_0]} \mathbf{1}_{\{\tau(B) > k\}} V^{1/2}(X_{k+1}, I_{k+1}) \right] \leq C_\varepsilon \sum_{k=0}^{\infty} \gamma'_{[n/c_0]} \gamma_{[n/c_0]} < \infty$$

by (18). With Lebesgue's dominated convergence theorem, we deduce that

$$\mathbb{E} \left[\sum_{k=0}^{\infty} \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} \left| U_{k+1}^{(3)} \right| \right] \leq C_\varepsilon \sum_{k=0}^{\infty} \gamma'_{[n/c_0]} \gamma_{[n/c_0]} < \infty$$

which implies that $\sum_{k=0}^n \gamma'_{a_k} \mathbf{1}_{\{\tau(B) > k\}} U_{k+1}^{(3)}$ converges almost surely to a finite limit. This completes the proof of the lemma. \square

We are now in position to prove Theorem 4.1. We start with (i).

Proposition 6.1. *Assume (A1-3) and let $B > 0$ be given. Then $|\theta_n - \theta^*| \mathbf{1}_{\{\tau(B) > n\}} \rightarrow 0$ with probability one as $n \rightarrow \infty$.*

Proof. The idea of the proof is borrowed from [11]. We saw in (22) that

$$\theta_{n+1} = \theta_n + \gamma_{a_n} H(\theta_n, I_{n+1}) + \gamma_{a_n}^2 r_n,$$

where $H = (H_1, \dots, H_d)$, $r_n = (r_{n,1}, \dots, r_{n,d})$, $H_i(\theta, I) = \theta(i) \left(\mathbf{1}_{\{I=i\}} - \theta(I) \right)$ and $r_{i,n} = -\theta(i)\theta(I) \frac{\mathbf{1}_{\{I=i\}} - \theta(I)}{1 + \gamma_{a_n} \theta(I)}$. We note that $|H| \leq 1$ and $|r_n| \leq 1$. Let $\varepsilon \in (0, \varepsilon^*)$ be such that when $\tau(B) \leq \xi(\varepsilon)$ (Lemma 6.5). We recall that the mean field function of the recursion is $h_i(\theta) = (\theta^*(i) - \theta(i)) / \sum_{j=1}^d \frac{\theta^*(j)}{\theta(j)}$. We introduce $\theta'_n = \theta_n + \sum_{j=n}^{\infty} \gamma_{a_j} \mathbf{1}_{\{\tau(B) > j\}} [H(\theta_j, I_{j+1}) - h(\theta_j)]$. From Lemma 6.6, $|\theta'_n - \theta_n| \rightarrow 0$ almost surely. $\{\theta'_n\}$ satisfies the recursion

$$\theta'_{n+1} = \theta'_n + \gamma_{a_n} h(\theta_n) + \gamma_{a_n} \mathbf{1}_{\{\tau(B) \leq n\}} [H(\theta_n, I_{n+1}) - h(\theta_n)] + \gamma_{a_n}^2 r_n.$$

We can then deduce:

$$\begin{aligned} |\theta'_{n+1} - \theta^*|^2 \mathbf{1}_{\{\tau(B) > n\}} &= |\theta'_n - \theta^*|^2 \mathbf{1}_{\{\tau(B) > n\}} + 2\gamma_{a_n} \langle \theta'_n - \theta^*, h(\theta_n) \rangle \mathbf{1}_{\{\tau(B) > n\}} + \gamma_{a_n} \mathbf{1}_{\{\tau(B) > n\}} r'_n \\ &\leq (1 - 2\varepsilon\gamma_{a_n}) |\theta'_n - \theta^*|^2 \mathbf{1}_{\{\tau(B) > n\}} + \gamma_{a_n} \mathbf{1}_{\{\tau(B) > n\}} (r'_n + \langle \theta'_n - \theta^*, h(\theta_n) - h(\theta'_n) \rangle). \end{aligned}$$

where $r'_n \rightarrow 0$ almost surely as $n \rightarrow \infty$. Since θ_n remains in the compact Θ_ε , h is continuous and since $|\theta'_n - \theta_n| \rightarrow 0$, it follows that $\langle \theta'_n - \theta^*, h(\theta_n) - h(\theta'_n) \rangle \rightarrow 0$. We can summarize the situation like this. Writing $U_n = |\theta'_n - \theta^*|^2 \mathbf{1}_{\{\tau(B) > n\}}$, we have:

$$U_{n+1} \leq (1 - 2\varepsilon\gamma_{a_n})U_n + \gamma_{a_n}r''_n, \quad (25)$$

where $r''_n \rightarrow 0$ as $n \rightarrow \infty$. This implies that $U_n \rightarrow 0$ which, given Lemma 6.6 proves the Proposition.

To see why $U_n \rightarrow 0$, let $\delta > 0$ be given. Take $n_0 > 0$ such that for $n \geq n_0$, $|r''_n| \leq 2\varepsilon\delta$ and $(1 - 2\varepsilon\gamma_{a_n})\mathbf{1}_{\{\tau(B) > n\}} > 0$. Then for $n \geq n_0$, $(U_{n+1} - \delta) \leq (1 - 2\varepsilon\gamma_{a_n})(U_n - \delta) + 2\varepsilon\gamma_{a_n}(r''_n/2\varepsilon - \delta) \leq (1 - 2\varepsilon\gamma_{a_n})(U_n - \delta)$ which implies that $\limsup(U_n - \delta) \leq 0$ and since $\delta > 0$ is arbitrary, we conclude that $\lim U_n = 0$. □

Proposition 6.2. *Assume (A1-3) and let $B > 0$ be given. For any function $f \in L_{V^{1/2}}$, denoting $\bar{f} = f - \pi^*(f)$, we have:*

$$\frac{1}{n} \sum_{k=1}^n \bar{f}(X_k, I_k) \mathbf{1}_{\{\tau(B) > k-1\}} \rightarrow 0, \text{ a.s. as } n \rightarrow \infty. \quad (26)$$

Proof. In view of Proposition 6.1 and (20), we only need to show that

$$\frac{1}{n} \sum_{k=1}^n \left(f(X_k, I_k) - \pi_{\theta_{k-1}}(f) \right) \mathbf{1}_{\{\tau(B) > k-1\}} \rightarrow 0 \text{ a.s.} \quad (27)$$

Kronecker's Lemma applied to (24) of Lemma 6.6 with $\gamma'_n = 1/n$, $H_\theta = f$ yields (27). □

6.3. Proof of Theorem 4.2

Proof. Assume that [a] hold. Define $\alpha_k = (1 + \gamma_0)^{(1+n_0)k}$ and suppose that $\limsup_{n \rightarrow \infty} n(v_n(i) - v_n(j)) = \infty$. This implies the existence of an increasing sequence of integers $\{n_k, k \geq 1\}$ such that $n_k(v_{n_k}(i) - v_{n_k}(j)) > \alpha_k$ and $(X_{n_k}, I_{n_k}) \in \mathcal{X}_i \times \{i\}$ but $(X_{n_k+n_0}, I_{n_k+n_0}) \notin \mathcal{X}_j \times \{j\}$ for all $k \geq 1$. Clearly, $n_k(v_{n_k}(i) - v_{n_k}(j)) > \alpha_k$ implies that $\theta_{n_k}(i)/\theta_{n_k}(j)$ converges to $+\infty$. But, since i leads to j , we can then find $\varepsilon > 0$ and k_0 such that for $k \geq k_0$:

$$\Pr[(X_{n_k+n_0}, I_{n_k+n_0}) \notin \mathcal{X}_j \times \{j\} | \mathcal{F}_{n_k}, (X_{n_k}, I_{n_k}) \in \mathcal{X}_i \times \{i\}, n_k(v_{n_k}(i) - v_{n_k}(j)) > \alpha_k] \leq (1 - \varepsilon).$$

Thus $\Pr(\limsup_{n \rightarrow \infty} n(v_n(i) - v_n(j)) = \infty) \leq \lim_{k \rightarrow \infty} (1 - \varepsilon)^k = 0$.

Assume that **[b]** hold. Define $\alpha_k = (1 + \gamma_0)^{2k}$ and suppose that $\limsup_{n \rightarrow \infty} n(\max_i v_n(i) - \min_j v_n(j)) = \infty$. Then we can find $i_0 \in \{1, \dots, d\}$ and an increasing sequence of integers $\{n_k, k \geq 1\}$ such that $\min_j v_n(j) = v_n(i_0)$, $n_k(\max_j v_{n_k}(j) - v_{n_k}(i_0)) > \alpha_k$, $(X_{n_k}, I_{n_k}) \notin \mathcal{X}_{i_0} \times \{i_0\}$ and $(X_{n_k+1}, I_{n_k+1}) \notin \mathcal{X}_{i_0} \times \{i_0\}$ for all $k \geq 1$. Then we can proceed as above and conclude. \square

Acknowledgments: The authors are grateful to Christophe Andrieu, David P. Landau, Eric Moulines and for very helpful discussions. We also thank David P. Sanders and the referees for comments that have made this paper hopefully more readable. This work is partly supported by the National Science Foundation grant DMS 0244638 and by a postdoctoral fellowship from the Natural Sciences and Engineering Research Council of Canada.

References

- [1] ANDRIEU, C. and ATCHADE, Y. F. (2007). On the efficiency of adaptive mcmc algorithms. *Electronic Communications in Probability* **12** 336–349.
- [2] ANDRIEU, C. and MOULINES, É. (2006). On the ergodicity properties of some adaptive MCMC algorithms. *Ann. Appl. Probab.* **16** 1462–1505.
- [3] ANDRIEU, C., MOULINES, É. and PRIOURET, P. (2005). Stability of stochastic approximation under verifiable conditions. *SIAM J. Control Optim.* **44** 283–312 (electronic).
- [4] ANDRIEU, C. and ROBERT, C. P. (2001). Controlled mcmc for optimal sampling. *Technical report, Université Paris Dauphine, Ceremade 0125* .
- [5] ATCHADE, Y. F. and LIU, J. S. (2006). Discussion of the paper by kou, zhou and wong. *Annals of Statistics* **To appear**.
- [6] ATCHADE, Y. F. and ROSENTHAL, J. S. (2005). On adaptive markov chain monte carlo algorithm. *Bernoulli* **11** 815–828.
- [7] BAXENDALE, P. H. (2005). Renewal theory and computable convergence rates for geometrically ergodic markov chains. *Annals of Applied Probability* **15** 700–738.
- [8] BENVENISTE, A., MÉTIVIER, M. and PRIOURET, P. (1990). *Adaptive Algorithms and Stochastic approximations*. Applications of Mathematics, Springer, Paris-New York.
- [9] BERG, B. A. and NEUHAUS, T. (1992). Multicanonical ensemble: A new approach to simulate first-order phase transitions. *Phys. Rev. Lett.* **68**.

- [10] BROCKWELL, A. E. and KADANE, J. B. (2005). Identification of regeneration times in MCMC simulation, with application to adaptive schemes. *J. Comput. Graph. Statist.* **14** 436–458.
- [11] DELYON, B. (1996). General results on the convergence of stochastic algorithms. *IEEE Trans. Automat. Control* **41** 1245–1255.
- [12] GEYER, C. J. and THOMPSON, E. (1995). Annealing markov chain monte carlo with applications to pedigree analysis. *Journal of the American Statistical Association* **90** 909–920.
- [13] GILKS, W. R., ROBERTS, G. O. and SAHU, S. K. (1998). Adaptive Markov chain Monte Carlo through regeneration. *J. Amer. Statist. Assoc.* **93** 1045–1054.
- [14] GREEN, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82** 711–732.
- [15] HAARIO, H., SAKSMAN, E. and TAMMINEN, J. (2001). An adaptive metropolis algorithm. *Bernoulli* **7** 223–242.
- [16] LIANG, F., LIU, C. and CARROLL, R. J. (2007). Stochastic approximation in monte carlo computation. *JASA* **102** 305–320.
- [17] LIANG, F. and WONG, W. H. (2001). Real-parameter evolutionary monte carlo with applications to bayesian mixture models. *Journal of the American Statistical Association* **96** 653–666.
- [18] MARINARI, E. and PARISI, G. (1992). Simulated tempering: A new monte carlo schemes. *Europhysics letters* **19** 451–458.
- [19] MIRA, A. and SARGENT, D. J. (2003). A new strategy for speeding Markov chain Monte Carlo algorithms. *Stat. Methods Appl.* **12** 49–60.
- [20] ROSENTHAL, J. S. and ROBERTS, G. O. (2007). Coupling and ergodicity of adaptive mcmc. *Journal of Applied Probability* **44** 458–475.
- [21] WANG, F. and LANDAU, D. P. (2001). Efficient, multiple-range random walk algorithm to calculate the density of states. *Physical Review Letters* **86** 2050–2053.