

COMPLEXITY BOUNDS FOR MARKOV CHAIN MONTE CARLO ALGORITHMS VIA DIFFUSION LIMITS

GARETH O. ROBERTS,* *University of Warwick*
JEFFREY S. ROSENTHAL,** *University of Toronto*

Abstract

We connect known results about diffusion limits of Markov chain Monte Carlo (MCMC) algorithms to the computer science notion of algorithm complexity. Our main result states that any weak limit of a Markov process implies a corresponding complexity bound (in an appropriate metric). We then combine this result with previously-known MCMC diffusion limit results to prove that under appropriate assumptions, the random-walk Metropolis algorithm in d dimensions takes $O(d)$ iterations to converge to stationarity, while the Metropolis-adjusted Langevin algorithm takes $O(d^{1/3})$ iterations to converge to stationarity.

Keywords: MCMC; convergence; complexity; diffusion limit; random-walk Metropolis algorithm; Metropolis-adjusted Langevin algorithm

2010 Mathematics Subject Classification: Primary 60J05; 60J25
Secondary 62F10; 62F15

1. Introduction

In the computer science literature, algorithms are often analysed in terms of ‘complexity’ bounds. In the Markov chain Monte Carlo (MCMC) literature, algorithms are sometimes understood in terms of diffusion limits. The purpose of this paper is to connect these two approaches, and in particular to show that diffusion limits (and other weak limits) can imply algorithm complexity bounds.

Complexity results in computer science go back at least to Cobham (1965), and took on greater focus with the pioneering *NP-complete* work of Cook (1971). In the Markov chain context, computer scientists have been bounding convergence times of Markov chain algorithms since at least Jerrum and Sinclair (1989), focusing largely on spectral gap bounds for Markov chains on finite state spaces. More recently, attention has turned to bounding spectral gaps of modern Markov chain algorithms on general (e.g. uncountable) state spaces, again primarily via spectral gaps; see, e.g. Woodard *et al.* (2009a), (2009b). These bounds often focus on the order of the convergence time in terms of some particular parameter, such as the dimension d of the corresponding state space.

Meanwhile, in statistics, MCMC algorithms are extremely widely used and studied (see, e.g. Brooks *et al.* (2011), and the many references therein), and their running times are an extremely important practical issue. They have been studied from a variety of perspectives, including

Received 19 August 2014; revision received 12 May 2015.

* Postal address: Department of Statistics, University of Warwick, Coventry CV4 7AL, UK.

Email address: gareth.o.roberts@warwick.ac.uk.

Supported in part by EPSRC grants EP/20620/01 and EP/S61577/01.

** Postal address: Department of Statistics, University of Toronto, Toronto, Ontario, M5S 3G3, Canada.

Email address: jeff@math.toronto.edu.

Supported in part by NSERC of Canada.

directly bounding the convergence in total variation distance (see, e.g. Rosenthal (1995a), (1996), (2002), Jones and Hobert (2001), (2004), and the references therein), convergence ‘diagnostics’ via statistical analysis of the Markov chain output (e.g. Gelman and Rubin (1992)), and most notably by proving weak convergence limits of sped up versions of the algorithms to diffusion limits (e.g. Roberts *et al.* (1997), Roberts and Rosenthal (1998)).

Direct total variation bounds for MCMC are sometimes presented in terms of the convergence order; see, e.g. Rosenthal (1995b) for order bounds for a Gibbs sampler for a variance components model. In addition, the MCMC diffusion limits often involve speeding up the original algorithm by a certain order, and then proving weak convergence to a fixed process which converges in $O(1)$ iterations, thus giving them the flavour of complexity order bounds too. However, these MCMC results are typically not stated precisely in terms of convergence time complexity results, and (perhaps because of this) they are often overlooked by the computer science complexity community.

In this paper we attempt to connect these two streams of Markov chain convergence time bounds. In particular, we establish (Theorem 1) that results about weak limits do directly imply corresponding complexity bounds (using an appropriate convergence metric as described below). We then apply our theorem to previous results about diffusion limits of MCMC algorithms (Section 3), to establish running time complexity order bounds for such MCMC algorithms as the random-walk Metropolis algorithm (Theorem 2) and the Metropolis-adjusted Langevin algorithm (Theorem 3).

2. Assumptions and main result

In this section we state our general result about obtaining convergence complexity bounds from weak limits. To set it up, let $(\mathcal{X}, \mathcal{F}, \rho)$ be a general measurable metric space, i.e. a nonempty (and possibly uncountable) set \mathcal{X} endowed with a metric ρ which induces a Borel σ -algebra \mathcal{F} of measurable subsets. We wish to bound the convergence of a stochastic process $\{X_t\}$ on $(\mathcal{X}, \mathcal{F})$ to its stationary probability distribution π . To measure the distance to stationarity, on finite state spaces one often (see, e.g. Aldous and Fill (2014, Section 2.4.1)) uses the total variation distance defined by

$$\|\mathcal{L}_x(X_t) - \pi\|_{\text{TV}} := \sup_{|f| \leq 1} |\mathbb{E}_x[f(X_t)] - \pi(f)|,$$

where the supremum is taken over all measurable functions $f: \mathcal{X} \rightarrow \mathbb{R}$ with $|f(x)| \leq 1$ for all $x \in \mathcal{X}$. Here $\mathcal{L}_x(X_t)$ is the law of X_t conditional on starting at $X_0 = x$, $\mathbb{E}_x[f(X_t)]$ is the expected value of f with respect to this law, and $\pi(f) = \int f(x)\pi(dx)$ is the expected value of f with respect to π .

This total variation distance can also be used on general state spaces in many instances; see, e.g. Rosenthal (1996). However, it is not appropriate for bounding the weak convergence which arises in the diffusion context, since it may not go to 0 for processes which converge only weakly to stationarity, so we do not use it here. Instead, we let

$$\text{Lip}_1^1 = \{f: \mathcal{X} \rightarrow \mathbb{R}, |f(x) - f(y)| \leq \rho(x, y) \text{ for all } x, y \in \mathcal{X}, |f| \leq 1\}$$

be the set of all functions from \mathcal{X} to \mathbb{R} with Lipschitz constant ≤ 1 and with $|f(x)| \leq 1$ for all $x \in \mathcal{X}$, and use the distance function

$$\|\mathcal{L}_x(X_t) - \pi\|_{\text{KR}} := \sup_{f \in \text{Lip}_1^1} |\mathbb{E}_x[f(X_t)] - \pi(f)|.$$

(Here ‘KR’ stands for ‘Kantorovich–Rubinstein’; see the proof of Proposition 1 below.) The distance $\|\cdot\|_{\text{KR}}$ is similar to, but more restrictive than, the total variation distance, and as discussed below (see Proposition 1); it metrises weak convergence and so is appropriate for our purposes.

We also note that many approaches to stationary instead directly bound the spectral gap of the corresponding Markov operator (e.g. Woodard *et al.* (2009a)). However, on general state spaces, the spectral gap is zero for Markov chains which are not ‘geometrically ergodic’ (see, e.g. Roberts and Rosenthal (1997, Theorem 2)), even if they do converge to stationarity. Furthermore, many MCMC algorithms are not geometrically ergodic (e.g. the random-walk Metropolis algorithm on target distributions with heavier-than-exponential tails; see Mengersen and Tweedie (1996, Theorem 3.3)). They also are often not reversible, which makes spectral gaps harder to study or interpret. For these reasons, we do not wish to restrict attention to spectral gaps, which is another reason that we use the metric $\|\cdot\|_{\text{KR}}$.

A related issue is what initial states X_0 should be considered. On finite state spaces, one often (e.g. Jerrum and Sinclair (1989, Section 2)) considers the worst case, by taking a supremum over all initial states x , i.e. using something like $\sup_{x \in \mathcal{X}} \|\mathcal{L}_x(X_t) - \pi\|_{\text{TV}}$. But this supremum is also frequently inappropriate on general state spaces. For instance, if \mathcal{X} is unbounded then as t increases one can start from worse and worse states X_0 so that the supremum might never go to 0. Instead, we need to specify more precisely which initial state(s) X_0 to consider. As a concrete choice, we will take the π -average of the distances to stationarity from all initial states X_0 in \mathcal{X} . That is, for any Markov chain $\{X_t\}$ on $(\mathcal{X}, \mathcal{F})$ with stationary distribution π , we measure the distance to stationarity at time t by the distance function

$$\mathbb{E}_{X_0 \sim \pi} \|\mathcal{L}_{X_0}(X_t) - \pi\|_{\text{KR}} := \int_{x \in \mathcal{X}} \pi(dx) \|\mathcal{L}_x(X_t) - \pi\|_{\text{KR}}.$$

Using this distance function, we can state our main result concerning bounding convergence to stationarity using weak convergence of a sequence of processes to another fixed process. To avoid technicalities, we assume that this limiting process is *càdlàg*, i.e. has sample paths which are continuous on the right with limits on the left. (In our MCMC examples, the limiting process will in fact be a diffusion with *continuous* sample paths, so this is not a problem.)

Theorem 1. *Let $X^{(d)} = \{X_t^{(d)}\}_{t \geq 0}$ be a stochastic process on $(\mathcal{X}, \mathcal{F}, \rho)$, for each $d \in \mathbb{N}$, which converges weakly in the Skorokhod topology as $d \rightarrow \infty$ to a *càdlàg* process $X^{(\infty)} = \{X_t^{(\infty)}\}_{t \geq 0}$, i.e. $X^{(d)} \xrightarrow{w} X^{(\infty)}$. Assume these processes all have the same stationary probability distribution π , and that $X^{(\infty)}$ converges (either weakly or in total variation distance) to π . Then, for any $\varepsilon > 0$, there are $D < \infty$ and $T < \infty$ such that*

$$\mathbb{E}_{X_0^{(d)} \sim \pi} \|\mathcal{L}_{X_0^{(d)}}(X_t^{(d)}) - \pi\|_{\text{KR}} < \varepsilon, \quad t \geq T, d \geq D.$$

Theorem 1 may be summarised as saying that if a sequence $\{X^{(d)}\}$ of Markov processes converges weakly to a limiting ergodic process, then we can bound the convergence of the sequence of processes uniformly over all sufficiently large d , i.e. the processes converge in $O(1)$ iterations with respect to d . We will next apply this result to previously known diffusion limits of common MCMC algorithms.

3. Application to MCMC

Our primary interest is in the use of Theorem 1 to bound the complexity of MCMC algorithms. We begin with the most popular MCMC algorithm, the random-walk Metropolis

(RWM) algorithm. This algorithm proceeds, given a positive target probability density π_d on the state space \mathbb{R}^d , by running a Markov chain $\{\mathbf{Z}_n^d\}_{n=0}^\infty$ as follows. Given the value \mathbf{Z}_n^d , a proposed new state $\mathbf{Y}_{n+1}^d \sim \text{MVN}(\mathbf{Z}_n^d, \sigma_d^2)$ is chosen from a multivariate normal distribution centered at \mathbf{Z}_n^d , and then with probability $\min[1, \pi_d(\mathbf{Y}_{n+1}^d)/\pi_d(\mathbf{Z}_n^d)]$ the proposal is accepted and $\mathbf{Z}_{n+1}^d = \mathbf{Y}_{n+1}^d$, otherwise with the remaining probability the proposal is rejected and $\mathbf{Z}_{n+1}^d = \mathbf{Z}_n^d$. This algorithm is easily seen to be irreducible and aperiodic and to leave π_d stationary, so it will converge asymptotically to π_d . The question then becomes how quickly it will converge, and what choice of proposal variance σ_d^2 is optimal.

In this context, Roberts *et al.* (1997) proved the result that $U^d \xrightarrow{w} U$ as $d \rightarrow \infty$, where $U_t^d = \mathbf{Z}_{\lfloor dt \rfloor, 1}^d$ is the first coordinate of the RWM algorithm sped up by a factor of d , U is a limiting ergodic Langevin diffusion, and \xrightarrow{w} indicates weak convergence in the usual Skorokhod topology. They proved this result under certain strong technical assumptions, namely that π_d takes on the special product form $\pi_d(\mathbf{x}) = \prod_{i=1}^d h(x_i)$ for some fixed function $h: \mathbb{R} \rightarrow (0, \infty)$ with h'/h Lipschitz continuous, and

$$\int \left[\frac{h'(x)}{h(x)} \right]^8 h(x) dx < \infty \quad (1)$$

and

$$\int \left[\frac{h''(x)}{h(x)} \right]^4 h(x) dx < \infty. \quad (2)$$

They also assumed the processes \mathbf{Z}^d are in stationarity, and that $\sigma_d^2 = \ell^2/(d-1)$ for some fixed $\ell > 0$.

This theorem of Roberts *et al.* (1997) allowed them to study the limiting diffusion U as a function of the proposal variance parameter ℓ , and optimise it to prove that the algorithm converges fastest when its asymptotic acceptance rate is equal to 0.234... (see also Roberts and Rosenthal (2001)). Furthermore, since their process U^d involved speeding up the original algorithm by a factor of d , their results seemed to imply that RWM required $O(d)$ iterations to converge. However, a precise statement of such a complexity bound was not provided.

In light of Theorem 1 above, we are now able to use the diffusion limit of Roberts *et al.* (1997) to give an actual complexity bound on the RWM algorithm. We need one slight technical extension, namely to replace (1) and (2) above by the slightly stronger conditions

$$\int \left[\frac{h'(x)}{h(x)} \right]^{12} h(x) dx < \infty \quad (3)$$

and

$$\int \left[\frac{h''(x)}{h(x)} \right]^6 h(x) dx < \infty. \quad (4)$$

We then have the following result, proved in Section 5 below.

Theorem 2. *Let $Z^{(d)}$ be a RWM algorithm satisfying the above technical assumptions of Roberts *et al.* (1997), except with (1) and (2) replaced by (3) and (4). Then, for any $\varepsilon > 0$, there is $D < \infty$ and $T < \infty$ such that*

$$\mathbb{E}_{Z_{0,1}^{(d)} \sim h} \left\| \mathcal{L}_{Z_{0,1}^{(d)}}(Z_{\lfloor dt \rfloor, 1}^{(d)}) - h \right\|_{\text{KR}} < \varepsilon, \quad t \geq T, d \geq D.$$

(Here $\mathcal{L}_{Z_{0,1}^{(d)}}(Z_{\lfloor dt \rfloor, 1}^{(d)})$ represents the probability distribution of the first coordinate of $Z^{(d)}$ at iteration number equal to the greatest integer not exceeding dt , conditional on the process

starting with its first coordinate equal to the specified state $Z_{0,1}^{(d)}$, and with all other coordinates of $Z_0^{(d)}$ chosen independently according to the density h .) Hence, the RWM algorithm takes $O(d)$ iterations to converge to within ε of stationarity in its first (or any one) coordinate.

We believe this to be the first precise general result about the complexity order of the RWM algorithm. It does require strong technical assumptions, but it still applies to a fairly general collection of densities on \mathbb{R}^d . Furthermore, empirical studies (see, e.g. Roberts and Rosenthal (2001)) indicate that even when RWM algorithms do not satisfy the technical assumptions, they still exhibit similar limiting behaviour.

Another MCMC diffusion limit concerns the Metropolis-adjusted Langevin algorithm (MALA). This algorithm is similar to the above RWM algorithm, except that now the proposal state $Y_{n+1}^d \sim \text{MVN}(Z_n^d + \frac{1}{2}\sigma_d^2 \nabla \log \pi_d(Z_n^d), \sigma_d^2)$ is chosen from a multivariate normal distribution centered at $Z_n^d + \frac{1}{2}\sigma_d^2 \nabla \log \pi_d(Z_n^d)$ (to better approximate π_d), and the above acceptance probability is modified by the ratio of the corresponding proposal normal distributions. In this context, Roberts and Rosenthal (1998) proved that $U^d \xrightarrow{w} U$, where $U_t^d = Z_{\lfloor d^{1/3}t \rfloor, 1}^d$ is the first coordinate of the MALA algorithm sped up by a factor of $d^{1/3}$, and U is again a limiting ergodic Langevin diffusion. This result again requires strong technical assumptions, this time that $\pi_d(\mathbf{x}) = \prod_{i=1}^d h(x_i)$ for some fixed function $h: \mathbb{R} \rightarrow (0, \infty)$ with polynomially-bounded log-derivatives of all orders, and finite moments of all orders, with h'/h Lipschitz continuous. They also assume that the processes Z^d are in stationarity, and that $\sigma_d^2 = \ell^2 d^{-1/3}$ for some fixed $\ell > 0$.

Roberts and Rosenthal's (1998) theorem allowed them to optimise the limiting diffusion U as a function of ℓ , and to prove that the algorithm converges fastest when its asymptotic acceptance rate is equal to 0.574... Also, since their process U^d involved speeding up the original algorithm by a factor of $d^{1/3}$, their results seemed to imply that MALA required $O(d^{1/3})$ iterations to converge. Once again, we can use Theorem 1 above to obtain the following more formal complexity bound (proved in Section 5 below).

Theorem 3. *Let $Z^{(d)}$ be a MALA algorithm on a product density in d dimensions satisfying the above technical assumptions of Roberts and Rosenthal (1998). Then, for any $\varepsilon > 0$, there is $D < \infty$ and $T < \infty$ such that (with notation as in Theorem 2)*

$$\mathbb{E}_{Z_{0,1}^{(d)} \sim h} \|\mathcal{L}_{Z_{0,1}^{(d)}}(Z_{\lfloor d^{1/3}t \rfloor, 1}^{(d)}) - h\|_{\text{KR}} < \varepsilon, \quad t \geq T, d \geq D.$$

Hence, the MALA algorithm takes $O(d^{1/3})$ iterations to converge to within ε of stationarity in its first (or any one) coordinate.

Finally, we note that a number of other diffusion limits have been proven for MCMC algorithms in other contexts. For example, Bédard (2007), (2008) and Sherlock and Roberts (2009) have extended the original RWM diffusion limit to more general target distributions; Roberts (1998), Neal and Roberts (2006), (2008), (2011), and Jourdain *et al.* (2013), (2014) have extended it to other related cases; and Neal *et al.* (2012) have established diffusion limits for RWM algorithms on discontinuous target densities. Each of these diffusion limit results could also be combined with Theorem 1 above to yield complexity order bounds in new contexts.

4. Proof of Theorem 1

In this section we prove Theorem 1. Along the way, we establish that $\|\cdot\cdot\cdot\|_{\text{KR}}$ metrises weak convergence (Proposition 1), and that $\mathbb{E}_{X_0 \sim \pi} \|\mathcal{L}_{X_0}(X_t^{(d)}) - \pi\|_{\text{KR}}$ is a nonincreasing function of t (Lemma 4). We first establish that $\|\cdot\cdot\cdot\|_{\text{KR}}$ is a norm.

Lemma 1. *Let S be any nonempty collection of functionals $\mathcal{X} \rightarrow \mathbb{R}$ which is symmetric (i.e. if $f \in S$ then $-f \in S$). Let $\|\mu\| = \sup_{f \in S} \mu(f)$. Then $\|\cdot\|$ is a (possibly infinite) norm function on the set of all signed measures on $(\mathcal{X}, \mathcal{F})$. In particular, $\|\cdot\|_{\text{KR}}$ is a norm.*

Proof. It is immediate that $\|0\| = 0$, and that $\|a\mu\| = a\|\mu\|$ for $a > 0$. The symmetry of S implies that $\|-\mu\| = \|\mu\|$. Finally, for the triangle inequality, we check that

$$\|\mu + \nu\| = \sup_{f \in S} (\mu(f) + \nu(f)) \leq \left(\sup_{f \in S} \mu(f) \right) + \left(\sup_{f \in S} \nu(f) \right) = \|\mu\| + \|\nu\|.$$

Hence, $\|\cdot\|$ is a norm. The claim about $\|\cdot\|_{\text{KR}}$ then follows by taking $S = \text{Lip}_1^1$. \square

We next show that truncating the metric ρ does not change Lip_1^1 .

Lemma 2. *Let $\rho^* = \min(2, \rho)$. Then*

$$\text{Lip}_1^1 = \{f: \mathcal{X} \rightarrow \mathbb{R}, |f(x) - f(y)| \leq \rho^*(x, y) \text{ for all } x, y \in \mathcal{X}, |f| \leq 1\}.$$

Proof. This is immediate since we always have $|f(x) - f(y)| \leq 2$ for $f \in \text{Lip}_1^1$. \square

Proposition 1. *The metric $\Delta(\mu, \nu) := \|\mu - \nu\|_{\text{KR}}$ metrises weak convergence of probability measures on $(\mathcal{X}, \mathcal{F}, \rho)$. That is, if $\{\mu_t\}$ and μ are probability measures on $(\mathcal{X}, \mathcal{F}, \rho)$, then $\{\mu_t\} \xrightarrow{w} \mu$ if and only if $\lim_{t \rightarrow \infty} \Delta(\mu_t, \mu) = 0$.*

Proof. Let ρ^* be as in Lemma 2. We first note that since ρ and ρ^* agree for distances less than or equal to 2, they give rise to precisely the same open subsets. Therefore, (\mathcal{X}, ρ^*) induces the same Borel σ -algebra \mathcal{F} that (\mathcal{X}, ρ) does and, thus, gives rise to the same Skorokhod topology. Hence, weak convergence on $(\mathcal{X}, \mathcal{F}, \rho)$ is precisely equivalent to weak convergence on $(\mathcal{X}, \mathcal{F}, \rho^*)$. Furthermore, by Lemma 2, the metric $\|\cdot\|_{\text{KR}}$ is the same on $(\mathcal{X}, \mathcal{F}, \rho^*)$ as on $(\mathcal{X}, \mathcal{F}, \rho)$. Hence, it suffices to prove the result on the truncated space $(\mathcal{X}, \mathcal{F}, \rho^*)$.

Now, since $(\mathcal{X}, \mathcal{F}, \rho^*)$ is a bounded metric space, it is known (see, e.g. Givens and Shortt (1984, Proposition 4)) that weak convergence on $(\mathcal{X}, \mathcal{F}, \rho^*)$ is metrised by the Wasserstein metric W_1 on (\mathcal{X}, ρ^*) , defined by

$$W_1(\mu, \nu) := \inf \mathbb{E}[\rho^*(X, Y)],$$

where the infimum is taken over all pairs (X, Y) of random variables on $(\mathcal{X}, \mathcal{F})$ such that $\mathcal{L}(X) = \mu$ and $\mathcal{L}(Y) = \nu$. On the other hand, again since $(\mathcal{X}, \mathcal{F}, \rho^*)$ is a bounded metric space, it is known (Kantorovich and Rubinstein (1958); see, e.g. Givens and Shortt (1984, p. 233)) that for probability measures μ and ν on Wasserstein metric $W_1(\mu, \nu)$ is precisely equal to $\|\mu - \nu\|_{\text{KR}}$. Combining these two facts, the result follows for $(\mathcal{X}, \mathcal{F}, \rho^*)$ and, hence, also for $(\mathcal{X}, \mathcal{F}, \rho)$. \square

Lemma 3. *If $X^{(\infty)}$ converges to π , either weakly or in total variation distance, then for all $x \in \mathcal{X}$ and $\varepsilon > 0$ there is $T < \infty$ such that $\|\mathcal{L}_x(X_T^{(\infty)}) - \pi\|_{\text{KR}} \leq \varepsilon/2$ for all $t \geq T$.*

Proof. If the convergence is weak then this follows from Proposition 1. If the convergence is in total variation distance then this still follows since $\|\cdot\|_{\text{KR}} \leq \|\cdot\|_{\text{TV}}$. \square

We are now in a position to prove convergence of $X_T^{(d)}$ for certain fixed times T .

Proposition 2. *Under the assumptions of Theorem 1, for any $x \in \mathcal{X}$ and $\varepsilon > 0$, there is $D < \infty$ and $T < \infty$ such that*

$$\|\mathcal{L}_x(X_T^{(d)}) - \pi\|_{\text{KR}} < \varepsilon, \quad d \geq D.$$

Proof. Using Lemma 1, by the triangle inequality, we have

$$\|\mathcal{L}_x(X_t^{(d)}) - \pi\|_{\text{KR}} \leq \|\mathcal{L}_x(X_t^{(d)}) - \mathcal{L}_x(X_t^{(\infty)})\|_{\text{KR}} + \|\mathcal{L}_x(X_t^{(\infty)}) - \pi\|_{\text{KR}}. \quad (5)$$

To continue, we recall that since $X^{(d)}$ converges weakly to $X^{(\infty)}$, it follows that $X_t^{(d)}$ converges weakly to $X_t^{(\infty)}$ for all fixed times $t > 0$ such that $X^{(\infty)}$ has probability 0 of *jumping* at time t , i.e. for all but at most a countable number of times t (since $X^{(\infty)}$ is càdlàg). By Lemma 3, there is $T < \infty$ such that $\|\mathcal{L}_x(X_T^{(\infty)}) - \pi\|_{\text{KR}} \leq \varepsilon/2$, and by increasing T as necessary we can assume that $X^{(\infty)}$ has probability 0 of jumping at time T . Then $X_T^{(d)}$ converges weakly to $X_T^{(\infty)}$, so by Proposition 1 there is $D < \infty$ such that, for all $d \geq D$, $\|\mathcal{L}_x(X_T^{(d)}) - \mathcal{L}_x(X_T^{(\infty)})\|_{\text{KR}} < \varepsilon/2$. The result then follows from (5). \square

Remark 1. If the weak convergence of $X^{(d)}$ to $X^{(\infty)}$ is assumed to be uniform over bounded time intervals, then we can strengthen Proposition 2 to say that for any $x \in \mathcal{X}$ and $\varepsilon > 0$ and $S < \infty$, there are $D < \infty$ and $T < \infty$ such that $\|\mathcal{L}_x(X_t^{(d)}) - \pi\|_{\text{KR}} < \varepsilon$ for all $t \in [T, T+S]$.

Corollary 1. *Under the assumptions of Theorem 1, for any $\varepsilon > 0$, there is $D < \infty$ and $T < \infty$ such that*

$$\mathbb{E}_{X_0 \sim \pi} \|\mathcal{L}_{X_0}(X_T^{(d)}) - \pi\|_{\text{KR}} < \varepsilon, \quad d \geq D.$$

Proof. We first let

$$A_m = \left\{ x \in \mathcal{X} : \|\mathcal{L}_x(X_t^{(\infty)}) - \pi\|_{\text{KR}} < \frac{\varepsilon}{4} \text{ for all } t \geq m \right\}.$$

Then $A_{m+1} \supseteq A_m$ by inspection, and $\bigcup_m A_m = \mathcal{X}$ by Lemma 3. Hence, by continuity of probabilities (see, e.g. Rosenthal (2000, Proposition 3.3.1)), $\lim_{m \rightarrow \infty} \pi(A_m) = 1$. We can therefore find $T < \infty$ such that $\pi(A_T) \geq 1 - (\varepsilon/8)$. As in the proof of Proposition 2, by increasing T as necessary we can assume that $X^{(\infty)}$ has probability 0 of jumping at time T .

Next, for this fixed T , let

$$B_m = \left\{ x \in \mathcal{X} : \|\mathcal{L}_x(X_T^{(d)}) - \mathcal{L}_x(X_T^{(\infty)})\|_{\text{KR}} < \frac{\varepsilon}{4} \text{ for all } d \geq m \right\}.$$

Then $B_{m+1} \supseteq B_m$ by inspection, and $\bigcup_m B_m = \mathcal{X}$ since $X_T^{(d)} \xrightarrow{w} X_T^{(\infty)}$, so again by continuity of probabilities we can find $D \in \mathbb{N}$ such that $\pi(B_D) \geq 1 - (\varepsilon/8)$.

We then compute that for this fixed T and D , and, for any $d \geq D$,

$$\begin{aligned} \mathbb{E}_{X_0 \sim \pi} \|\mathcal{L}_{X_0}(X_T^{(d)}) - \pi\|_{\text{KR}} &= \mathbb{E}_{X_0 \sim \pi} (\mathbf{1}_{\{X_0 \in A_T \cap B_D\}} \|\mathcal{L}_{X_0}(X_T^{(d)}) - \pi\|_{\text{KR}} \\ &\quad + \mathbb{E}_{X_0 \sim \pi} (\mathbf{1}_{\{X_0 \notin A_T \cap B_D\}} \|\mathcal{L}_{X_0}(X_T^{(d)}) - \pi\|_{\text{KR}}) \\ &\leq \left[\left(\frac{\varepsilon}{4} \right) + \left(\frac{\varepsilon}{4} \right) \right] + \left[\left(\frac{\varepsilon}{8} \right) + \left(\frac{\varepsilon}{8} \right) \right] \times 2 \\ &= \varepsilon, \end{aligned}$$

where for the first term we have used the triangle inequality, and for the second term we have used the fact that by definition we always have $\|\mathcal{L}_x(X_T^{(d)}) - \pi\|_{\text{KR}} \leq 2$ for any x and d . This completes the proof. \square

Corollary 1 is nearly what we need to prove Theorem 1. However, for Theorem 1 we want the convergence to be within ε for *all* $t \geq T$, not just for one fixed T (nor just for all t in some bounded time interval, as in Remark 1). Unfortunately, $\|\mathcal{L}_x(X_t^{(d)})\|_{\text{KR}} - \pi$ might not be a nonincreasing function of t (though $\|\mathcal{L}_x(X_t^{(d)})\|_{\text{TV}} - \pi$ always is; see, e.g. Roberts and Rosenthal (2004, Proposition 3(c))). On the other hand, fortunately the quantity $\mathbb{E}_{X_0 \sim \pi} \|\mathcal{L}_{X_0}(X_t^{(d)})\|_{\text{KR}} - \pi$ is indeed nonincreasing.

Lemma 4. *Let $\|\cdot\|$ be any norm function on signed measures on $(\mathcal{X}, \mathcal{F})$. Let $P^t(x, \cdot)$ be the transition probabilities for a Markov chain on $(\mathcal{X}, \mathcal{F})$ with stationary probability distribution π . Let $\text{dist}(t) = \mathbb{E}_{X_0 \sim \pi} \|P^t(X_0, \cdot) - \pi\|$. Then $\text{dist}(t)$ is a nonincreasing function of t . In particular, in the context of Theorem 1, $\mathbb{E}_{X_0 \sim \pi} \|\mathcal{L}_{X_0}(X_t^{(d)})\|_{\text{KR}} - \pi$ is a nonincreasing function of t .*

Proof. We compute by stationarity that, for $s, t > 0$,

$$\begin{aligned} \text{dist}(s+t) &= \mathbb{E}_{X_0 \sim \pi} \|P^{s+t}(X_0, \cdot) - \pi\| \\ &= \mathbb{E}_{X_0 \sim \pi} \left\| \int_{y \in \mathcal{X}} P^s(X_0, dy) P^t(y, \cdot) - \pi \right\| \\ &\leq \mathbb{E}_{X_0 \sim \pi} \int_{y \in \mathcal{X}} P^s(X_0, dy) \|P^t(y, \cdot) - \pi\| \\ &= \mathbb{E}_{Y_0 \sim \pi} \|P^t(Y_0, \cdot) - \pi\| \\ &= \text{dist}(t); \end{aligned}$$

thus, proving the first claim. The claim about $\mathbb{E}_{x \sim \pi} \|\mathcal{L}_x(X_t^{(d)})\|_{\text{KR}} - \pi$ then follows by Lemma 1 upon setting $P^t(x, A) = \mathbb{P}[X_t^{(d)} \in A \mid X_0^{(d)} = x]$. \square

Theorem 1 then follows by combining Corollary 1 and Lemma 4.

5. Proofs of Theorems 2 and 3

Theorems 2 and 3 *nearly* follow immediately by applying Theorem 1 to the diffusion limit results of Roberts *et al.* (1997) and of Roberts and Rosenthal (1998), respectively. However, there is one technical issue. The previous diffusion limit results assume that the process begins in the stationary distribution. By contrast, Theorem 1 involves $\mathcal{L}_x(X_t)$, i.e. *conditioning* on the stochastic processes' first coordinate beginning at a specific state $x \in \mathcal{X}$. So, to prove Theorems 2 and 3, we need to establish that the diffusion limit results remain valid even upon conditioning on the starting value of the processes.

For the RWM algorithm, this does indeed follow, at least upon strengthening (1) and (2) to (3) and (4) as above.

Proposition 3. *Let $Z^{(d)}$ be a RWM algorithm satisfying the above technical assumptions of Roberts *et al.* (1997), except with (1) and (2) replaced by (3) and (4). Then, for π -almost everywhere $x \in \mathcal{X}$, ${}_x U^d \xrightarrow{w} {}_x U$ as $d \rightarrow \infty$, where ${}_x U_t^d = (\mathbf{Z}_{[dt],1}^d \mid \mathbf{Z}_{0,1}^d = x)$ is the first coordinate of the RWM algorithm sped up by a factor of d , conditional on starting at the state x , and ${}_x U$ is the limiting ergodic Langevin diffusion U also conditional on starting at x .*

Proof. The proof is very similar to the proof of the unconditioned diffusion limit theorem of Roberts *et al.* (1997), using generators and cores (cf. Ethier and Kurtz (1986)). The only point of departure between our proof and theirs concerns their Lemma 2.1, which states that for each

fixed $t > 0$,

$$\lim_{d \rightarrow \infty} \mathbb{P}[Z_s^{(d)} \in F_d \text{ for all } 0 \leq s \leq t] = 1,$$

where F_d is the event that both $|A_d| < d^{-1/8}$ and $|B_d| < d^{-1/8}$, where

$$A_d := \frac{1}{d-1} \sum_{i=2}^d (((\log h(x_i))')^2 - \mathbb{E}_{x_i \sim h}[(\log h(x_i))']^2)$$

and

$$B_d := \frac{1}{d-1} \sum_{i=2}^d ((\log h(x_i))'' - \mathbb{E}_{x_i \sim h}[(\log h(x_i))'']).$$

To complete our proof, we need to show that this statement remains valid, even when conditioning on starting at a specific state $x \in \mathcal{X}$. (In fact, it would suffice to replace $\frac{1}{8}$ by any other power $\alpha \in (0, \frac{1}{2})$, but we do not need to do that.)

To this end, fix $t > 0$, and let

$$p(d, x) = \mathbb{P}[Z_s^{(d)} \notin F_d \text{ for some } 0 \leq s \leq t \mid Z_{0,1}^{(d)} = x],$$

and let $r(d) = \mathbb{E}_{x \sim h} p(d, x)$ be its expected value when averaged over x in stationarity. Also, let

$$v_j = \mathbb{E}_{x_i \sim h} [((\log h(x_i))')^2 - \mathbb{E}_{x_i \sim h}[(\log h(x_i))']^2]^j$$

and recall that $|v_j| < \infty$ for $1 \leq j \leq 6$ by (3). Then, by expanding out the power $(A_d)^6$, while omitting terms involving v_1 since clearly $v_1 = 0$, we conclude that

$$\begin{aligned} \mathbb{E}_\pi[(A_d)^6] &= (d-1)^{-6} \left[\binom{d-1}{1} v_6 + (d-1)(d-2) \binom{6}{2} v_2 v_4 + \binom{d-1}{2} \binom{6}{3} v_3^2 \right. \\ &\quad \left. + \binom{d-1}{3} \binom{6}{2, 2, 2} v_2^3 \right] \\ &= (d-1)^{-6} [(d-1)v_6 + 15(d-1)(d-2)v_2 v_4 + 10(d-1)(d-2)v_3^2 \\ &\quad + 15(d-1)(d-2)(d-3)v_2^3]. \end{aligned}$$

In particular, $\mathbb{E}_\pi[(A_d)^6] = O(d^{-3})$ as $d \rightarrow \infty$. Hence, by Markov's inequality,

$$\mathbb{P}_\pi(|A_d| > d^{-1/8}) \leq \frac{\mathbb{E}[(A_d)^6]}{(d^{-1/8})^6} = O(d^{-3+(6/8)}) = O(d^{-9/4}).$$

Similarly, $\mathbb{P}_\pi(|B_d| > d^{-1/8}) \leq O(d^{-9/4})$, so that also $\mathbb{P}[Z_s^{(d)} \notin F_d] \leq O(d^{-9/4})$.

Next, we note that there are $O(dt)$ different RWM iterations corresponding to times s with $0 \leq s \leq t$. Hence, by subadditivity,

$$r(d) := \mathbb{P}_\pi[Z_s^{(d)} \notin F_d \text{ for some } 0 \leq s \leq t] \leq O(dt d^{-9/4}) = O(td^{-5/4}).$$

In particular, $\sum_{d=2}^{\infty} r(d) < \infty$, which is the key.

Finally, we wish to show that $\lim_{d \rightarrow \infty} p(d, x) = 0$ with probability 1. To that end, let $\varepsilon > 0$ and set $S_d = \{x \in \mathcal{X} : p(d, x) \geq \varepsilon\}$. Then writing $\mathbb{P}_h(A)$ for $\int_A h(x) dx$, it follows

by Markov's inequality that $\mathbb{P}_h(S_d) \leq \mathbb{E}_h[p(d, x)]/\varepsilon = r(d)/\varepsilon$. Hence, $\sum_{d=2}^{\infty} \mathbb{P}_h(S_d) \leq \sum_{d=2}^{\infty} r(d)/\varepsilon < \infty$. Hence, by the Borel–Cantelli lemma, $\mathbb{P}_h(S_d \text{ infinitely often}) = 0$, i.e. the set of x with an infinite sequence of d with $p(d, x) \geq \varepsilon$ has probability 0. This means that with probability 1, $\limsup_{d \rightarrow \infty} p(d, x) < \varepsilon$. Since this holds for all $\varepsilon > 0$, it follows that with probability 1, $\lim_{d \rightarrow \infty} p(d, x) = 0$, as desired; thus, completing the proof. \square

Theorem 2 then follows by combining Theorem 1 and Proposition 3.

Finally, we prove a similar result for the MALA algorithm. In this case, no strengthening of the assumptions is required.

Proposition 4. *Let $Z^{(d)}$ be a MALA algorithm on a product density in d dimensions satisfying the above technical assumptions of Roberts and Rosenthal (1998). Then, for π -almost everywhere $x \in \mathcal{X}$, ${}_x U^d \xrightarrow{w} {}_x U$ as $d \rightarrow \infty$, where ${}_x U_t^d = (\mathbf{Z}_{\lfloor d^{1/3}t \rfloor, 1}^d \mid \mathbf{Z}_{0,1}^d = x)$ is the first coordinate of the RWM algorithm sped up by a factor of $d^{1/3}$, conditional on starting at x , and ${}_x U$ is the limiting ergodic Langevin diffusion U also conditional on starting at x .*

Proof. The proof involves modifying the weak convergence proof of Roberts and Rosenthal (1998), along lines very similar to that of Proposition 3, so we omit the details. Furthermore, since Roberts and Rosenthal (1998) assumed finite moments of *all* polynomial orders, there is no need to strengthen any of their assumptions as was necessary for Proposition 3. \square

Theorem 3 then follows by combining Theorem 1 and Proposition 4.

Acknowledgements

We thank Dawn Woodard and Alexandre Thiéry for helpful discussions of these matters, and thank the anonymous referee for several very valuable comments.

References

- ALDOUS, D. AND FILL, J. A. (2014). Reversible Markov chains and random walks on graphs. Unfinished monograph. Available at <http://www.stat.berkeley.edu/~aldous/RWG/book.html>.
- BÉDARD, M. (2007). Weak convergence of Metropolis algorithms for non-i.i.d. target distributions. *Ann. Appl. Prob.* **17**, 1222–1244.
- BÉDARD, M. (2008). Optimal acceptance rates for Metropolis algorithms: Moving beyond 0.234. *Stoch. Process. Appl.* **118**, 2198–2222.
- BROOKS, S., GELMAN, A., JONES, G. L. AND MENG, X.-L. (eds) (2011). *Handbook of Markov chain Monte Carlo*. Chapman & Hall/CRC, Boca Raton, FL.
- COBHAM, A. (1965). The intrinsic computational difficulty of functions. In *Proceedings of the 1964 International Congress for Logic, Methodology, and Philosophy of Science*, North-Holland, Amsterdam, pp. 24–30.
- COOK, S. A. (1971). The complexity of theorem-proving procedures. In *Proc. Third Annual ACM Symposium on Theory of Computing*, ACM, New York, pp. 151–158.
- ETHIER, S. N. AND KURTZ, T. G. (1986). *Markov Processes: Characterization and Convergence*. John Wiley, New York.
- GELMAN, A. AND RUBIN, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statist. Sci.* **7**, 457–472.
- GIVENS, C. R. AND SHORTT, R. M. (1984). A class of Wasserstein metrics for probability distributions. *Michigan Math. J.* **31**, 231–240.
- JONES, G. L. AND HOBERT, J. P. (2001). Honest exploration of intractable probability distributions via Markov chain Monte Carlo. *Statist. Sci.* **16**, 312–334.
- JONES, G. L. AND HOBERT, J. P. (2004). Sufficient burn-in for Gibbs samplers for a hierarchical random effects model. *Ann. Statist.* **32**, 784–817.
- JOURDAIN, B., LELIÈVRE, T. AND MIAOJEDOW, B. (2015). Optimal scaling for the transient phase of the Metropolis Hastings algorithm: The mean-field limit. *Ann. Appl. Prob.* **25**, 2263–2300

- JOURDAIN, B., LELIÈVRE, T. AND MIASOJEDOW, B. (2014). Optimal scaling for the transient phase of Metropolis Hastings algorithms: the longtime behavior. *Bernoulli* **20**, 1930–1978.
- KANTOROVIČ, L. AND RUBINŠTEJN, G. Š. (1958). On a space of completely additive functions. *Vestnik Leningrad. Univ.* **13**, 52–59.
- MENGERSEN, K. L. AND TWEEDIE, R. L. (1996). Rates of convergence of the Hastings and Metropolis algorithms. *Ann. Statist.* **24**, 101–121.
- NEAL, P. AND ROBERTS, G. (2006). Optimal scaling for partially updating MCMC algorithms. *Ann. Appl. Prob.* **16**, 475–515.
- NEAL, P. AND ROBERTS, G. (2008). Optimal scaling for random walk Metropolis on spherically constrained target densities. *Methodol. Comput. Appl. Prob.* **10**, 277–297.
- NEAL, P. AND ROBERTS, G. (2011). Optimal scaling of random walk Metropolis algorithms with non-Gaussian proposals. *Methodol. Comput. Appl. Prob.* **13**, 583–601.
- NEAL, P., ROBERTS, G. AND YUEN, W. K. (2012). Optimal scaling of random walk Metropolis algorithms with discontinuous target densities. *Ann. Appl. Prob.* **22**, 1880–1927.
- ROBERTS, G. O. (1998). Optimal Metropolis algorithms for product measures on the vertices of a hypercube. *Stoch. Stoch. Reports* **62**, 275–283.
- ROBERTS, G. O. AND ROSENTHAL, J. S. (1997). Geometric ergodicity and hybrid Markov chains. *Electron. Commun. Prob.* **2**, 13–25.
- ROBERTS, G. O. AND ROSENTHAL, J. S. (1998). Optimal scaling of discrete approximations to Langevin diffusions. *J. R. Statist. Soc. B* **60**, 255–268.
- ROBERTS, G. O. AND ROSENTHAL, J. S. (2001). Optimal scaling for various Metropolis–Hastings algorithms. *Statist. Sci.* **16**, 351–367.
- ROBERTS, G. O. AND ROSENTHAL, J. S. (2004). General state space Markov chains and MCMC algorithms. *Prob. Surv.* **1**, 20–71.
- ROBERTS, G. O., GELMAN, A. AND GILKS, W. R. (1997). Weak convergence and optimal scaling of random walk Metropolis algorithms. *Ann. Appl. Prob.* **7**, 110–120.
- ROSENTHAL, J. S. (1995a). Minorization conditions and convergence rates for Markov chain Monte Carlo. *J. Amer. Statist. Assoc.* **90**, 558–566, 1136.
- ROSENTHAL, J. S. (1995b). Rates of convergence for Gibbs sampler for variance components models. *Ann. Statist.* **23**, 740–761.
- ROSENTHAL, J. S. (1996). Analysis of the Gibbs sampler for a model related to James–Stein estimators. *Statist. Comput.* **6**, 269–275.
- ROSENTHAL, J. S. (2000). *A First Look at Rigorous Probability Theory*. World Scientific, River Edge, NJ.
- ROSENTHAL, J. S. (2002). Quantitative convergence rates of Markov chains: a simple account. *Electron. Commun. Prob.* **7**, 123–128.
- SHERLOCK, C. AND ROBERTS, G. O. (2009). Optimal scaling of the random walk Metropolis on elliptically symmetric unimodal targets. *Bernoulli* **15**, 774–798.
- WOODARD, D. B., SCHMIDLER, S. C. AND HUBER, M. L. (2009a). Conditions for rapid mixing of parallel and simulated tempering on multimodal distributions. *Ann. Appl. Prob.* **19**, 617–640.
- WOODARD, D. B., SCHMIDLER, S. C. AND HUBER, M. L. (2009b). Sufficient conditions for torpid mixing of parallel and simulated tempering. *Electron. J. Prob.* **14**, 780–804.