

1 Hierarchical models and the tuning of random walk Metropolis
2 algorithms

3 Mylène Bédard
4 Département de mathématiques et de statistique
5 Université de Montréal
6 Montréal, Canada, H3C 3J7
7 bedard@dms.umontreal.ca

8 **Abstract**

We obtain weak convergence and optimal scaling results for the random walk Metropolis algorithm with a Gaussian proposal distribution. The sampler is applied to hierarchical target distributions, which form the building block of many Bayesian analyses. The global asymptotically optimal proposal variance derived may be computed as a function of the specific target distribution considered. We also introduce the concept of locally optimal tunings, *i.e.* tunings that depend on the current position of the Markov chain. The theorems are proved by studying the generator of the first and second components of the algorithm, and verifying their convergence to the generator of a modified RWM algorithm and a diffusion process, respectively. The rate at which the algorithm explores its state space is optimized by studying the speed measure of the limiting diffusion process. We illustrate the theory with two examples. Applications of these results on simulated and real data are also presented.

9 **AMS 2000 subject classifications:** Primary 60F05; secondary 65C40.

10 **Keywords:** acceptance probability; diffusion process; efficiency; position-dependent re-
11 versibility; RWM-within-Gibbs; weak convergence

12 1. Introduction

13 Random walk Metropolis (RWM) algorithms are widely used to sample from complex or
14 multidimensional probability distributions ([15], [12]). The simplicity and versatility of these
15 samplers often make them the default option in the MCMC toolbox. Implementing a RWM
16 algorithm involves a tuning step, to ensure that the process explores its state space as fast as
17 possible, and that the sample produced be representative of the probability distribution of
18 interest (the target distribution). In this paper, we solve an aspect of the tuning problem for
19 a large class of target distributions with correlated components. This issue has mainly been
20 studied for product target densities, but attention has recently turned towards more complex
21 target models ([7], [14]). The specific type of target distribution considered here is formed of
22 components which are related according to a hierarchical structure. These distributions are
23 ubiquitous in several fields of research (finance, biostatistics, physics, to name a few), and
24 constitute the basis of many Bayesian inferences.

25 Bayesian hierarchical models are comprised of a likelihood function $f(\mathbf{d}|\theta)$, which is the
26 statistical model for the observed data \mathbf{d} . The parameters θ are then modeled using a prior
27 distribution $\pi(\theta|\rho)$; since this prior might not be easy to determine, it is common practice to
28 assume that the hyperparameters ρ are themselves distributed according to a non-informative
29 prior distribution $\pi(\rho)$. The various models thus represent different levels of hierarchy and
30 give rise to a posterior distribution $\pi(\theta, \rho|\mathbf{d})$, which is often quite complex. Most of the time,
31 this distribution cannot be studied analytically or sampled directly, and thus simulation
32 algorithms such as MCMC methods are required to perform a statistical analysis. Samplers
33 such as the RWM, RWM-within-Gibbs, and Adaptive Metropolis (see [11]) are usually the
34 default algorithms for such targets.

35 The idea behind RWM algorithms is to build a Markov chain having the Bayesian posterior
36 (target) distribution as its stationary distribution. To implement this method, users must
37 select a proposal distribution from which are generated candidates for the Markov chain.
38 This distribution should ideally be similar to the target, while remaining accessible from a
39 sampling viewpoint. A pragmatic choice is to let the proposed moves be normally distributed
40 around the latest value of the sample. Tuning the variance of the normal proposal distri-
41 bution (σ^2) has a significant impact on the speed at which the sampler explores its state
42 space (hereafter referred to as “efficiency”), with extremal variances leading to slow-mixing
43 algorithms. In particular, large variances seldom induce suitable candidates and result in
44 lazy processes; small variances yield hyperactive processes whose tiny steps lead to a time-
45 consuming exploration of the state space. Seeking for an intermediate value that optimizes
46 the efficiency of the RWM algorithm, *i.e.* a proposal variance σ^2 offering sizable steps that
47 are still accepted a reasonable proportion of the time, is called the optimal scaling problem.

48 The optimal scaling issue of the RWM algorithm with a Gaussian proposal has been addressed
49 by many researchers over the last few decades. It has been determined in [17] that target
50 densities formed of independent and identically distributed (i.i.d.) components correspond
51 to an optimal proposal variance $\hat{\sigma}^2(n) \approx 5.66/\{n\mathbb{E}[(\log f(X))']\}$, where f is the density
52 of one target component and n the number of target components. This optimal proposal
53 variance has also been shown to correspond to an optimal expected acceptance rate of 23.4%,
54 where the acceptance rate is defined as the proportion of candidates that are accepted by
55 the algorithm. Generalizing this conclusion is an intricate task and further research on the

56 subject has mainly been restricted to the case of target distributions formed of independent
57 components (see [18], [16], [2], [3], [5], [6]). In the specific case of multivariate normal target
58 distributions however, the optimal variance and acceptance rate may be easily determined
59 (see [16], [1]). Lately, [7] and [14] have also performed scaling analyses of non-product target
60 densities. These advances are important, as MCMC methods are mainly used when dealing
61 with complex models, which only rarely satisfy the independence assumption among target
62 components. These results however assume that the correlation structure among target
63 components is known and used in generating candidates for the chain. This is a restrictive
64 assumption that leads, as expected, to an optimal acceptance rate of 23.4% (see [18] for an
65 explanation).

66 In this paper, we focus on solving the optimal scaling problem for a wide class of models that
67 include a dependence relationship, the hierarchical distributions. Weak convergence results
68 are derived without explicitly characterizing the dependency among target components, and
69 thus rely on a Gaussian proposal distribution with diagonal covariance matrix. The optimal
70 proposal variance may then be obtained from these results, *i.e.* by maximizing the speed
71 measure of the limiting diffusion process. This constitutes significant advances in under-
72 standing the theoretical underpinnings of the RWM sampler. More importantly in practice,
73 the results theoretically support the use of RWM-within-Gibbs over RWM samplers and provide
74 a convenient approach for obtaining a new type of proposal variances. These proposal
75 variances are a function of the current state of the Markov chain; they thus evolve with the
76 chain and lead to more appropriate candidates in the RWM-within-Gibbs algorithm.

77 In the next section, we describe the target distribution and introduce some notation related
78 to the RWM sampler. The theoretical optimal scaling results are stated in Section 3, and then
79 illustrated with two examples using RWM samplers in Section 4. In Section 5, the potential of
80 RWM-within-Gibbs with local scalings is illustrated in Bayesian contexts through a simulation
81 study and an application on real data. Extensions are briefly discussed in Section 6, while
82 appendices contain proofs.

83 2. Framework

84 Consider an n -dimensional target distribution consisting of a mixing component X_1 and of
85 $n-1$ conditionally i.i.d. components X_i ($i = 2, \dots, n$) given X_1 . Suppose that this distribution
86 has a target density π with respect to Lebesgue measure, where

$$\pi(\mathbf{x}) = f_1(x_1) \prod_{i=2}^n f(x_i | x_1) . \quad (1)$$

87 To obtain a sample from the target density in (1), we rely on a RWM algorithm with
88 a Gaussian proposal distribution. This sampler builds an n -dimensional Markov chain
89 $\{\mathbf{X}^{(n)}[j]; j \in \mathbb{N}\}$ having $\pi(\mathbf{x})$ as its stationary density. Given $\mathbf{X}^{(n)}[j] = \mathbf{x}$, the time- j state of
90 the Markov chain, one iteration is performed according to the following steps:

- 91 1. generate a candidate $\mathbf{Y}^{(n)}[j+1] = \mathbf{y}$ from a $\mathcal{N}(\mathbf{x}, D_n)$, where D_n is a diagonal variance
92 matrix with elements $(\sigma_1^2(n), \sigma^2(n), \dots, \sigma^2(n))$. In particular, set $D_n = \ell^2 I_n/n$, where
93 $\ell > 0$ is a tuning parameter and I_n the n -dimensional identity matrix;

- 94 2. compute the acceptance probability $\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ 1, \frac{\pi(\mathbf{y})}{\pi(\mathbf{x})} \right\}$;
- 95 3. generate $U[j + 1] \sim \mathcal{U}(0, 1)$;
- 96 4. if $U[j + 1] \leq \alpha(\mathbf{x}, \mathbf{y})$, accept the candidate and set $\mathbf{X}^{(n)}[j + 1] = \mathbf{y}$; otherwise, the
- 97 Markov chain remains at the same state for another time interval and $\mathbf{X}^{(n)}[j + 1] = \mathbf{x}$.

98 Optimal scaling results widely rely on the use of Gaussian proposal distributions which, due

99 to their symmetry, lead to a simplified form of the acceptance probability. Although generally

100 not emphasized in the literature, we note that the proposal variance could also be a function

101 of \mathbf{x} , which would result in a non-homogeneous random walk sampler. In that case, there

102 would be no simplification in the Metropolis-Hastings acceptance probability and Step 2

103 would then be replaced by

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ 1, \frac{\pi(\mathbf{y})q_n(\mathbf{x};\mathbf{y})}{\pi(\mathbf{x})q_n(\mathbf{y};\mathbf{x})} \right\} ,$$

104 where $q_n(\mathbf{y}; \mathbf{x})$ is the density of a $\mathcal{N}(\mathbf{x}, D_n(\mathbf{x}))$.

105 In what follows we work towards finding the optimal value of ℓ , *i.e.* leading to an optimally

106 mixing chain. The proofs of the theoretical results rely on CLTs and LLNs; as such, the

107 results are obtained by letting $n \rightarrow \infty$. This is a common approach in MCMC theory

108 and does not prevent users from applying the asymptotically optimal value of ℓ in lower

109 dimensional contexts (as small as $n = 10$ or 15). Indeed, a particularity of optimal scaling

110 results is that the asymptotic behaviour kicks in extremely rapidly, as shall be witnessed in

111 the examples of Section 4.

112 The first thought of most MCMC users when facing a target density as in (1) would be to use

113 a RWM-within-Gibbs algorithm, which consecutively updates subgroups of the n components

114 in a given iteration. The tuning of RWM-within-Gibbs algorithms has been addressed in [16],

115 but only for target distributions with i.i.d. components and Gaussian targets with correlation.

116 Focusing on RWM algorithms is thus a good starting point to understand the behaviour of

117 samplers applied to hierarchical target distributions. The results expounded in this paper

118 lead to the concept of local tunings, which is particularly appealing in the context of RWM-

119 within-Gibbs. Incidentally, the proofs in appendices provide a theoretical justification for the

120 use of locally optimal scalings in RWM-within-Gibbs, see [4]. These findings are illustrated

121 in the examples of Section 5.

122 In Sections 2.1, 2.2, and 3, we expound how to obtain asymptotically optimal variances

123 D_n and $D_n(\mathbf{x})$ for RWM and RWM-within-Gibbs, respectively. Section 2.1 describes the

124 regularity conditions imposed on $\pi(\mathbf{x})$, while Section 2.2 explains why the proposal matrix

125 $D_n = \ell^2 I_n/n$ is the optimal choice for obtaining the theoretical results that shall be presented

126 in Section 3.

127 2.1. Assumptions on the target density

128 To characterize the asymptotic behaviour of the conditionally i.i.d. components X_i ($i =$

129 $2, \dots, n$), we impose some regularity conditions on the densities f_1 and f in (1). The density

130 f_1 is assumed to be a continuous function on \mathbb{R} , with $\mathcal{X}_1 = \{x_1 : f_1(x_1) > 0\}$ forming an

131 open interval.

132 For all fixed $x_1 \in \mathcal{X}_1$, $f(x|x_1)$ is a positive \mathcal{C}^2 density on \mathbb{R} and $\frac{\partial}{\partial x} \log f(x|x_1)$ is Lipschitz
133 continuous with constant $K(x_1)$ such that $\mathbb{E}[K^2(X_1)] < \infty$. Here, \mathcal{C}^2 denotes the space of
134 real-valued functions with continuous second derivative. For all fixed $x \in \mathcal{X} = \mathbb{R}$, $f(x|x_1)$ is
135 a \mathcal{C}^2 function on \mathcal{X}_1 and $\frac{\partial}{\partial x} \log f(x|x_1)$ is Lipschitz continuous with constant $L(x)$ such that
136 $\mathbb{E}[L^4(X)] < \infty$. Furthermore,

$$\mathbb{E}_X \left[\left(\frac{\partial}{\partial X} \log f(X|x_1) \right)^4 \right] < \infty \quad \forall x_1 \in \mathcal{X}_1 \quad \text{with} \quad \mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \right] < \infty ; \quad (2)$$

137 hereafter, the notation $\mathbb{E}_X[\cdot]$ means that the expectation is computed with respect to X
138 conditionally on the other variables in the expression; the first expectation in (2) is thus
139 obtained according to the conditional distribution of X given X_1 . Where there is no confusion
140 possible, $\mathbb{E}[\cdot]$ shall be used to denote an expectation with respect to all random variables in
141 the expression. The above regularity conditions constitute an extension of those stated in
142 [3] for target distributions with independent components, and are weaker than would be a
143 Lipschitz continuity assumption on the bivariate function $\frac{\partial}{\partial x} \log f(x|x_1)$. They also imply
144 that the Lipschitz constants $K(x_1)$ and $L(x)$ themselves satisfy a Lipschitz condition.

145 We now impose further conditions on $f(x|x_1)$ to account for the movements of the coordinate
146 X_1 when studying the asymptotic behaviour of a component X_i ($i = 2, \dots, n$). These move-
147 ments should not be too abrupt so for almost all fixed $x \in \mathcal{X}$, $\frac{\partial}{\partial x_1} \log f(x|x_1)$ is Lipschitz
148 continuous with constant $M(x)$ such that $\mathbb{E}[M^2(X)] < \infty$ and

$$\mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right] < \infty \quad \forall x_1 \in \mathcal{X}_1 \quad \text{with} \quad \mathbb{E} \left[\left(\frac{\partial}{\partial X_1} \log f(X|X_1) \right)^2 \right] < \infty . \quad (3)$$

149 Finally, in order to characterize the asymptotic behaviour of the mixing component X_1 , we
150 introduce assumptions that are closely related to the Bernstein von Mises Theorem. Let
151 $\mathbf{X}_{2:n} = (X_2, \dots, X_n)$, $\mathbf{X} = (X_2, X_3, \dots)$, and \rightarrow_p denote convergence in probability. Assume
152 that $\mathbb{V}(X_1|\mathbf{X}_{2:n}) \rightarrow_p 0$, and denote $\mu \equiv \mu(\mathbf{X})$ such that $\mu_n \equiv \mu_n(\mathbf{X}_{2:n}) = \mathbb{E}[X_1|\mathbf{X}_{2:n}] \rightarrow_p \mu$
153 as $n \rightarrow \infty$, with $|\mu| < \infty$. Hereafter, we make a small abuse of notation by letting μ
154 and μ_n sometimes denote the random variable or the realisation, depending on the context.
155 Furthermore, define $\tilde{X}_1 = \sqrt{n}(X_1 - \mu_n)$; for almost all $\mathbf{x}_{2:n} \in \mathbb{R}^{n-1}$, the conditional density
156 of \tilde{X}_1 given $\mathbf{x}_{2:n}$, $f_1(\mu_n + \tilde{x}_1/\sqrt{n}|\mathbf{x}_{2:n})/\sqrt{n}$, is assumed to converge almost surely to $g_1(\tilde{x}_1|\mathbf{x})$,
157 a continuous density on \mathbb{R} with respect to Lebesgue measure. In fact, the information on X_1
158 increases linearly in n , meaning that the limiting density of $X_1|\mathbf{x}_{2:n}$ is degenerate, but that
159 a standard rescaling leads to a non-trivial density on \mathbb{R} (normal distribution).

160 2.2. Form of the proposal variance matrix D_n

161 In Section 3, we focus on deriving weak convergence and optimal scaling results for the RWM
162 algorithm with a Gaussian proposal by letting n , the dimension of the target density in
163 (1), approach ∞ . Traditionally, asymptotically optimal scaling results have been obtained
164 by studying the limiting path of a given component (X_2 say) as $n \rightarrow \infty$. In the case of
165 target distributions with i.i.d. components (and some extensions), the components of the
166 RWM algorithm are asymptotically independent of each other and their limiting behaviour
167 is regimented by identical one-dimensional Markovian processes. In the current correlated
168 framework, we expect the presence of an asymptotic dependence relationship among X_i

169 ($i \in \{2, \dots, n\}$) and X_1 , in the spirit of (1). In the following section, we thus study the
 170 limiting behaviour of components X_1 and X_2 separately, on their respective conditional space.
 171 This approach allows us to quantify the mixing rate of each component X_i conditionally on
 172 the others, and to propose optimal scalings for the sampler.

173 To obtain non-trivial limiting processes describing the behaviour of the RWM sampler as
 174 $n \rightarrow \infty$, we need to fix the form of the proposal scalings $\sigma_1^2(n), \sigma^2(n)$. Whilst the proposals
 175 are independent, a single accept-reject step is used, which makes the paths of the components
 176 dependent. We aim to choose the maximal scalings that avoid a degenerate limit (of either
 177 0 or 1) for this acceptance probability. Since the distribution of X_1 conditional on $\mathbf{X}_{2:n}$
 178 contracts at a rate of \sqrt{n} , then if $\sigma_1(n)/\sqrt{n} \rightarrow \infty$ the proposed jumps in X_1 will be too
 179 large. If $\sigma_1(n)/\sqrt{n} \rightarrow 0$, then the change in X_1 makes no contribution to the acceptance
 180 probability in the limit; to maximise movements we, therefore, require $\sigma_1(n) \propto 1/\sqrt{n}$. Now,
 181 the conditional distribution of $\mathbf{X}_{2:n}$ given X_1 does not contract with n . Nonetheless, when
 182 proposing jumps in $\mathbf{X}_{2:n}$ using $\sigma^2(n) = \sigma^2$, the odds of rejecting an n -dimensional candidate
 183 increase with n and lead to a degenerate (null) acceptance probability. To overcome this
 184 problem we then let the proposal variance be a decreasing function of the dimension. In fact,
 185 since Lipschitz conditions control the contribution to the accept-reject ratio coming from the
 186 movements of X_1 , a similar argument to that which leads to $\sigma(n) \propto 1/\sqrt{n}$ in the case of i.i.d.
 187 targets applies again here. We therefore set $D_n = \ell^2 I_n/n$, where $\ell > 0$ is a tuning parameter
 188 and I_n the n -dimensional identity matrix.

189 As $n \rightarrow \infty$, it becomes necessary to speed up time to compensate for the reduced movement
 190 along components $\mathbf{X}_{2:n}$. The time interval between each proposed candidate is thus set to
 191 $1/n$ and we study the continuous-time, sped up version of the initial Markov chain defined
 192 as $\{\mathbf{W}^{(n)}(t); t \geq 0\} = \{\mathbf{X}^{(n)}[\lfloor nt \rfloor]; t \geq 0\}$, where $\lfloor \cdot \rfloor$ is the floor function. Similarly to the
 193 i.i.d. case, a limiting diffusion is obtained for the rescaled one-dimensional process related to
 194 X_i ($i \geq 2$), but this time its behaviour is conditional on X_1 .

195 Since the first coordinate X_1 converges to a point μ , a transformation $\tilde{X}_1 = \sqrt{n}(X_1 - \mu_n)$ is
 196 required to obtain the limiting behaviour of this component. We thus study the continuous-
 197 time process $\{\tilde{\mathbf{W}}^{(n)}(t); t \geq 0\} = \{(\tilde{X}_1^{(n)}[\lfloor t \rfloor], \mathbf{X}_{2:n}^{(n)}[\lfloor t \rfloor]); t \geq 0\}$; in other words, we are now
 198 looking at a magnified, centered version of the path associated to X_1 . This transformation
 199 leads to proposal distributions $\tilde{Y}_1 = \sqrt{n}(Y_1 - \mu_n) \sim \mathcal{N}(\tilde{x}_1, \ell^2)$ and $Y_i \sim \mathcal{N}(x_i, \ell^2/n)$, $i =$
 200 $2, \dots, n$ with $\ell > 0$; it thus cancels the effect of n in $\sigma_1^2(n)$. Without the speed up of time,
 201 the limiting process for \tilde{X}_1 is then a propose-accept-reject on the conditional density for \tilde{X}_1 ,
 202 given the current values of $\mathbf{X}_{2:n}$; this is made precise in Theorem 1. When considering the
 203 diffusion limit for X_i ($i \geq 2$) with time sped-up, this effectively means that at every instant,
 204 X_1 is simply a sample from its conditional distribution given the current values of $\mathbf{X}_{2:n}$; this
 205 is made precise in Theorem 2.

206 We note that an alternative scaling of $\sigma_1(n) \propto 1/n$ could also be applied. The sped-up
 207 limiting process would then be a diffusion for all coordinates, and would be easier to describe.
 208 However, this would also be a deliberate handicapping of the algorithm since the change in
 209 X_1 would make no contribution to the acceptance probability in the limit. A suboptimal
 210 $\sigma_1^2(n)$, besides altering the movements of X_1 , would thus also indirectly affect the efficiency
 211 according to which $\mathbf{X}_{2:n}$ explores its state space.

212 **3. Asymptotics of the RWM algorithm**

213 In this section we introduce results about the limiting behaviour (as $n \rightarrow \infty$) of the time- and
 214 scale-adjusted univariate processes $\{\tilde{W}_1^{(n)}(t); t \geq 0\}$ and $\{W_i^{(n)}(t); t \geq 0\}$ ($i = 2, \dots, n$). From
 215 these results we determine the asymptotically optimal scaling (AOS) values and acceptance
 216 rate (AOAR) that optimize the mixing of the algorithm.

217 Hereafter, we let \Rightarrow denote weak convergence in the Skorokhod topology and $B(t)$ a Brownian
 218 motion at time t ; the cumulative distribution function of a standard normal random variable
 219 is denoted by $\Phi(\cdot)$.

220 **Theorem 1.** *Consider a RWM algorithm with proposal distribution $\mathcal{N}(\mathbf{x}, \ell^2 I_n/n)$ used to
 221 sample from a target density π as in (1). Suppose that π satisfies the conditions on f_1 and f
 222 specified in Section 2.1, and that $\mathbf{X}^{(n)}(0)$ is distributed according to π in (1).*

223 If $\frac{1}{n} \sum_{i=2}^n \left(\frac{\partial}{\partial X_i} \log f(X_i | X_1 = \mu_n + \frac{\tilde{X}_1}{\sqrt{n}}) \right)^2 \rightarrow_p \tilde{\gamma}(\mu)$ with

$$\tilde{\gamma}(\mu) = \mathbb{E}_{\mathbf{X}} \left[\left(\frac{\partial}{\partial X} \log f(X | \mu(\mathbf{X})) \right)^2 \right] = \int_{\mathbb{R}} \left(\frac{\partial}{\partial x} \log f(x | \mu(\mathbf{X})) \right)^2 f(x | \mu(\mathbf{X})) dx < \infty,$$

224 then the magnified process $\{\tilde{W}_1^{(n)}(t); t \geq 0\} \Rightarrow \{\tilde{W}_1(t); t \geq 0\}$. Here, $W_1(0)$ and $W_i(0)$
 225 ($i = 2, 3, \dots$) are distributed according to the densities f_1 and f respectively, which implies that
 226 $\tilde{W}_1(0)$ is distributed according to the density g_1 in Section 2.1. Given the time- t state $\tilde{\mathbf{W}}(t) =$
 227 $(\tilde{x}_1, \underline{\mathbf{x}})$, the process $\{\tilde{W}_1(t); t > 0\}$ evolves as the continuous-time version of a special RWM
 228 algorithm applied to the target density $g_1(\tilde{x}_1 | \underline{\mathbf{x}})$; the proposal distribution of this algorithm is
 229 a $\mathcal{N}(\tilde{x}_1, \ell^2)$ and the acceptance rule is defined as

$$\alpha^*(\tilde{x}_1, \tilde{y}_1 | \underline{\mathbf{x}}) = \Phi \left(\frac{\log \frac{g_1(\tilde{y}_1 | \underline{\mathbf{x}})}{g_1(\tilde{x}_1 | \underline{\mathbf{x}})} - \frac{\ell^2}{2} \tilde{\gamma}(\mu)}{\ell \tilde{\gamma}^{1/2}(\mu)} \right) + \frac{g_1(\tilde{y}_1 | \underline{\mathbf{x}})}{g_1(\tilde{x}_1 | \underline{\mathbf{x}})} \Phi \left(\frac{-\log \frac{g_1(\tilde{y}_1 | \underline{\mathbf{x}})}{g_1(\tilde{x}_1 | \underline{\mathbf{x}})} - \frac{\ell^2}{2} \tilde{\gamma}(\mu)}{\ell \tilde{\gamma}^{1/2}(\mu)} \right). \quad (4)$$

230 *Proof.* See Appendix A.1. □

231 This result describes the limiting path associated to the coordinate \tilde{X}_1 as $n \rightarrow \infty$, which is
 232 Markovian with respect to the history of the multidimensional chain $\mathcal{F}^{\tilde{\mathbf{W}}}(t)$. We recall that
 233 the conditional distribution of X_1 given $\mathbf{X}_{2:n}$ contracts at a rate of \sqrt{n} and that $\sigma_1(n) \propto 1/\sqrt{n}$.
 234 Conditionally on $\underline{\mathbf{X}}$, the transformed \tilde{X}_1 thus mixes according to $\mathcal{O}(1)$ and explores its
 235 conditional state space much more efficiently than the other components, as shall be witnessed
 236 in Theorem 2. The asymptotic process found can be described as an atypical one-dimensional
 237 RWM algorithm, whose acceptance rule $\alpha^*(\tilde{x}_1, \tilde{y}_1 | \underline{\mathbf{x}})$ and target density $g_1(\tilde{x}_1 | \underline{\mathbf{x}})$ both vary
 238 according to $\underline{\mathbf{x}}$ at every iteration. The acceptance function α^* in (4) satisfies the reversibility
 239 condition with respect to $g_1(\tilde{x}_1 | \underline{\mathbf{x}})$ (see [3] for more details about this acceptance function).

240 Theorem 1 is interesting from a theoretical perspective, but cannot be used to optimize
 241 the global mixing of the algorithm. Although we could try to determine the value of ℓ
 242 leading to the optimal mixing of X_1 on its conditional space, it will be wiser to focus instead
 243 on optimizing the mixing rate of $\mathbf{X}_{2:n}$ on its own conditional space given X_1 . Since the
 244 distribution of X_1 contracts about μ_n , the position of this coordinate heavily depends on the

245 current state of $\mathbf{X}_{2:n}$. We shall also see in Theorem 2 that given X_1 , the coordinates X_i
 246 ($i \geq 2$) explore their conditional state space according to $\mathcal{O}(n)$. Since these coordinates take
 247 more time exploring their conditional distribution and heavily affect the position of X_1 , then
 248 the global performance of the sampler is subjected to the mixing of $\mathbf{X}_{2:n}$ conditionally on
 249 X_1 .

250 **Theorem 2.** *Consider a RWM algorithm with proposal distribution $\mathcal{N}(\mathbf{x}, \ell^2 I_n/n)$ used to
 251 sample from a target density π as in (1). Suppose that π satisfies the conditions on f_1 and f
 252 specified in Section 2.1, and that $\mathbf{X}^{(n)}(0)$ is distributed according to π in (1).*

253 *For $i = 2, \dots, n$, we have $\{W_i^{(n)}(t); t \geq 0\} \Rightarrow \{W_i(t); t \geq 0\}$, where $W_i(0)$ ($i \geq 2$) is
 254 distributed according to f , and $W_1(0)$ according to f_1 . Conditionally on $W_1(t)$, the evolution
 255 of $\{W_i(t); t > 0\}$ over an infinitesimal interval dt satisfies*

$$dW_i(t) = v^{1/2}(\ell, W_1(t))dB(t) + \frac{1}{2}v(\ell, W_1(t))\frac{\partial}{\partial W_i(t)} \log f(W_i(t)|W_1(t)) dt, \quad (5)$$

256 *with*

$$v(\ell, x_1) = 2\ell^2 \mathbb{E}_{Z_1} \left[\Phi \left(-\frac{\ell}{2} \gamma^{1/2}(x_1, Z_1) \right) \right], \quad (6)$$

257 $Z_1 = \sqrt{n}(Y_1 - x_1)/\ell \sim \mathcal{N}(0, 1)$, and

$$\gamma(x_1, z_1) = z_1^2 \mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right] + \mathbb{E}_X \left[\left(\frac{\partial}{\partial X} \log f(X|x_1) \right)^2 \right]. \quad (7)$$

258 *Proof.* See Appendix A.2. □

259 Equation (5) describes the behaviour of the process at the next instant, $(t + dt)$, given its
 260 position at t . This expression should not come as a surprise: each rescaled component X_i
 261 ($i = 2, \dots, n$) asymptotically behaves according to a diffusion process that is Markovian
 262 with respect to $\mathcal{F}^{(W_1, W_i)}(t)$. Examination of (5) also tells us that $f(W_i(t)|W_1(t))$ is invariant
 263 for this diffusion process (see [19], for instance). We finally recall that $\sigma(n) \propto 1/\sqrt{n}$ and
 264 therefore, conditionally on X_1 , the rescaled X_i mixes according to $\mathcal{O}(n)$. Each coordinate
 265 X_i thus requires more iterations than were required by the coordinate X_1 to explore its
 266 conditional state space.

267 Since X_1 and X_i ($i \geq 2$) use different time rescaling factors, the asymptotic behaviour of these
 268 coordinates cannot be expressed as a bivariate diffusion process. To obtain such a diffusion,
 269 we would have to rely on inhomogeneous proposal variances to ensure that X_1 also mixes in
 270 $\mathcal{O}(n)$ iterations; as mentioned at the end of Section 2, this would require setting $\sigma_1(n) = \ell/n$,
 271 $\sigma(n) = \ell/\sqrt{n}$ for $\ell > 0$. This framework would of course be suboptimal as it would restrain
 272 the X_1 movements. Proposed jumps for X_1 would then become insignificant, and so the first
 273 term in (7) would be null.

274 **Remark 3.** *Studying the limiting behaviour of X_1 and X_i ($i = 2, \dots, n$) separately does
 275 not cause information loss. In fact, studying the paths of X_1, X_2 simultaneously would re-
 276 quire letting the test function h of the generator in (A.3) be a function of (X_1, X_2) . Such
 277 a generator would however be developed as an expression in which cross-derivative terms
 278 (e.g. $\frac{\partial^2}{\partial x_1 \partial x_2} h(x_1, x_2)$) are null, which confirms that given the current state of the asymptotic
 279 process, one-dimensional moves are performed independently for each coordinate.*

280 The limiting processes in Theorems 1 and 2 indicate that the component X_1 explores its
 281 conditional state space at a different (higher) rate than $\mathbf{X}_{2:n}$ explores its own. Combined to
 282 the specific Markovian forms of the limiting processes obtained (with respect to $\mathcal{F}^{\tilde{\mathbf{W}}}(t)$ and
 283 $\mathcal{F}^{(W_1, W_i)}(t)$ respectively), this points towards the need for updating X_1 and $\mathbf{X}_{2:n}$ separately,
 284 assessing the superiority of RWM-within-Gibbs samplers for sampling from hierarchical tar-
 285 gets. These algorithms update blocks of components successively, a design that allows fully
 286 exploiting the characteristics of the target considered. To our knowledge, this is the first time
 287 that asymptotic results are used to theoretically validate the superiority of RWM-within-
 288 Gibbs over RWM samplers for hierarchical target distributions. This theoretical superiority
 289 is obviously tempered in practice by an increased computational effort; the extent of this
 290 computational overhead is however difficult to quantify in full generality. To this end, Sec-
 291 tion 5 presents two examples that illustrate the performance of the RWM-within-Gibbs and
 292 compare it to RWM and Adaptive Metropolis samplers.

293 3.1. Optimal tuning of the RWM algorithm

294 To be confident that the n -dimensional chain has entirely explored its state space, we must
 295 be certain that every one-dimensional path has explored its own space. In the correlated
 296 framework considered, the overall mixing rate of the RWM sampler is only as fast as the
 297 slowest component. As explained in Section 3, optimal mixing of the algorithm shall be
 298 attained by optimizing the mixing of the coordinates X_i , $i = 2, \dots, n$. In the limit, the only
 299 quantity that depends on the proposal variance (*i.e.* on ℓ) is $v(\ell, W_1(t))$ in (6). To optimize
 300 mixing, it thus suffices to find the diffusion process that goes the fastest, *i.e.* the value of ℓ
 301 for which the speed measure $v(\ell, W_1(t))$ is optimized.

302 The speed measure in (6) is quite intuitive; it is in fact similar to that obtained when studying
 303 i.i.d. target densities. The main difference lies in the form of $\gamma(x, z)$ which, in the i.i.d. case, is
 304 given by the constant term $\gamma = \mathbb{E}[(\frac{\partial}{\partial X} \log f(X))^2]$. The second term in (7) is thus equivalent
 305 to γ , and consists in a measure of roughness of the conditional density $f(x_i|x_1)$ under a
 306 variation of x_i ($i \geq 2$). In the case of hierarchical target distributions, we find an extra term
 307 that might be viewed as a measure of roughness of $f(x_i|x_1)$ under a variation of x_1 . This
 308 term is weighted by z_1^2 , the square of the (standardized) candidate increment for the first
 309 component; in other words, the further the candidate y_1 is from the current x_1 , the greater
 310 is the weight attributed to the associated measure of roughness. Of course, in optimizing the
 311 speed measure function, we do not need to know in advance the exact value of the proposed
 312 standardized increment z_1 ; the speed measure averages over this quantity.

313 It is interesting to note that optimizing the speed measure leads to local proposal variances of
 314 the form $\hat{\ell}^2(W_1(t))/n$. Such proposal variances would then be used for proposing a candidate
 315 at the next instant $t + dt$, given the position of the mixing coordinate at time t . These local
 316 proposal variances thus vary from one iteration to another, by opposition to usual tunings in
 317 the literature that are fixed for the duration of the algorithm. Naturally, if both expectations
 318 in (7) are constant with respect to x_1 , then the proposal variance obtained by maximizing
 319 the speed measure also is constant.

320 **Remark 4.** *It turns out that local proposal variances optimizing (6) are bounded above by*
 321 $2.38/\mathbb{E}_X^{1/2}[(\frac{\partial}{\partial X} \log f(X|x_1))^2]$, *the asymptotically optimal scaling (AOS) values for targets*

322 with i.i.d. components given a fixed $X_1 = x_1$. Indeed, if $X_1 = x_1$ is fixed across itera-
 323 tions, we find ourselves in an i.i.d. setting and the associated speed measure is expressed as
 324 $2\ell^2\Phi(-\ell\mathbb{E}_X^{1/2}[(\frac{\partial}{\partial X}\log f(X|x_1))^2]/2)$. The mentioned upper bounds then follow from the fact
 325 that the function $\Phi(\cdot)$ in (6) decreases faster in ℓ than $\Phi(\cdot)$ in the above expression.

326 Relying on a local variance $\hat{\ell}(x_1)$ to propose a candidate for the next time interval is usually
 327 time-consuming, as it involves numerically solving for the appropriate local proposal variance
 328 at every iteration. Since the process is assumed to start in stationarity and X_1 explores its
 329 conditional state space faster than the other coordinates, we might determine a value $\hat{\ell}$ that
 330 is fixed across iterations by integrating the speed measure $v(\ell, \cdot)$ over \mathcal{X}_1 with respect to the
 331 marginal distribution f_1 . Hence, the global (unconditional) asymptotically optimal scaling
 332 value $\hat{\ell}$ maximizes the function

$$\begin{aligned}\mathbb{E}_{X_1}[v(\ell, X_1)] &= 2\ell^2\mathbb{E}_{X_1, Z_1}\left[\Phi\left(-\frac{\ell}{2}\gamma^{1/2}(X_1, Z_1)\right)\right] \\ &= 2\ell^2\int_{\mathcal{X}_1}\int_{\mathbb{R}}\Phi\left(-\frac{\ell}{2}\gamma^{1/2}(x_1, z_1)\right)\phi(z_1)f_1(x_1)dz_1dx_1,\end{aligned}$$

333 where $\phi(\cdot)$ is the probability density function of a standard normal random variable.

334 **Remark 5.** *The asymptotic process introduced in Theorem 2 naturally leads us to the concept*
 335 *of local proposal variances. It is however unclear whether the local tunings obtained by maxi-*
 336 *mizing (6) really optimize the mixing rate of the algorithm. Indeed, the proof of Theorem 2 is*
 337 *carried out with ℓ^2 constant; this allows, among other things, relying on the simplified form*
 338 *for the acceptance probability. In order to claim that the local proposal variances obtained are*
 339 *optimal, a weak convergence result would need to be proven using a general proposal variance*
 340 *of the form $\sigma^2(n, x_1) = \ell^2(x_1)/n$. This extension is not trivial, as the ratio of proposal den-*
 341 *sities $q_n(\mathbf{x}; \mathbf{y})/q_n(\mathbf{y}; \mathbf{x})$ would then need to be included in the acceptance probability. Since*
 342 *the concept of locally optimal proposal variances is numerically demanding in the current*
 343 *framework, we choose to focus on ℓ^2 constant.*

344 *In RWM-within-Gibbs, the blocks X_1 and $\mathbf{X}_{2:n}$ are updated consecutively and the situation is*
 345 *therefore different. In that case, local variances of the form $\sigma^2(n, x_1) = \ell^2(x_1)/n$ obtained by*
 346 *maximizing (6) may be used to update the block $\mathbf{X}_{2:n}$. Since X_1 is updated separately, the*
 347 *first term in (7) is null, which makes local variances easier to compute. Furthermore, since*
 348 *local variances only depend on X_1 (which is updated separately), the ratio $q_n(\mathbf{x}; \mathbf{y})/q_n(\mathbf{y}; \mathbf{x})$*
 349 *is equal to 1 and does not need to be included in the acceptance probability. Local variances*
 350 *are thus very appealing in that context and shall be studied in Section 5.*

351 Rather than tuning the sampler using the global AOS value, one may instead monitor the
 352 acceptance rate in order to work with an optimally mixing version of the RWM algorithm.
 353 To express optimal scaling results in terms of acceptance rates, we introduce the expected
 354 acceptance rate of the n -dimensional stationary RWM algorithm with a normal proposal:

$$a_n(\ell) = \int \int \alpha(\mathbf{x}, \mathbf{y}) \left(\frac{\ell}{\sqrt{n}}\right)^{-n} \phi_n\left(\frac{\mathbf{y} - \mathbf{x}}{\ell/\sqrt{n}}\right) \pi(\mathbf{x}) d\mathbf{y} d\mathbf{x},$$

355 where $\phi_n(\cdot)$ denotes the probability density function of an n -dimensional standard normal
 356 random variable. Optimal mixing results for the RWM sampler are summarized in the fol-
 357 lowing corollary.

358 **Corollary 6.** *In the settings of Theorem 2, the global asymptotically optimal scaling value $\hat{\ell}$*
 359 *maximizes*

$$2\ell^2 \int_{\mathcal{X}_1} \int_{\mathbb{R}} \Phi \left(-\frac{\ell}{2} \gamma^{1/2}(x, z) \right) \phi(z) f_1(x) \, dz \, dx .$$

360 *Furthermore, we have that*

$$\lim_{n \rightarrow \infty} a_n(\ell) = a(\ell) \equiv 2 \int_{\mathcal{X}_1} \int_{\mathbb{R}} \Phi \left(-\frac{\ell}{2} \gamma^{1/2}(x, z) \right) \phi(z) f_1(x) \, dz \, dx ,$$

361 *and the corresponding asymptotically optimal acceptance rate is given by $a(\hat{\ell})$.*

362 In contrast to the i.i.d. case, the AOAR found is not independent of the densities f_1 and f .
 363 Hence, there is not a huge advantage in choosing to tune the acceptance rate of the algorithm
 364 over the proposal variance; in fact, both approaches involve the same effort. Although it would
 365 also be possible to compute an overall acceptance rate associated to using local proposal
 366 variances, it could not be used to tune the algorithm. Building an optimal Markov chain
 367 based on local proposal variances would imply modifying the proposal variance at every
 368 iteration, which cannot be achieved by solely monitoring the acceptance rate.

369 For simplicity, the theoretical results expounded in this section attribute the same tuning
 370 constant ℓ to all n components. In practice, when a RWM algorithm is used to sample from
 371 a hierarchical target, users will likely want to use a different proposal variance for the mixing
 372 component X_1 . In fact, the proofs of Theorems 1 and 2 easily generalize to the case of
 373 inhomogeneous proposal variances.

374 **Corollary 7.** *Let $Y_1 \sim \mathcal{N}(x_1, \ell^2 \kappa_1^2/n)$ with $0 < \kappa_1 < \infty$ and $\mathbf{Y}_{2:n} \sim \mathcal{N}(\mathbf{x}_{2:n}, \ell^2 I_{n-1}/n)$,*
 375 *where $Y_1, \mathbf{Y}_{2:n}$ are independent. Then, Theorems 1 and 2 hold as stated, except that the*
 376 *limiting proposal distribution in Theorem 1 is $\tilde{Y}_1 \sim \mathcal{N}(\tilde{x}_1, \ell^2 \kappa_1^2)$ and the random variable Z_1*
 377 *in Theorem 2 is such that $Z_1 \sim \mathcal{N}(0, \kappa_1^2)$.*

378 In this paper, we consider the simple, yet useful hierarchical model described in (1) and
 379 featuring a single mixing component X_1 . This is a natural starting point to study weak
 380 convergence of RWM algorithms for hierarchical targets, and even for correlated targets in
 381 general. There exist many generalizations of (1), just as there are many extensions of the
 382 proposal distribution considered. Some extensions of the hierarchical target are considered
 383 in the discussion, but we do not aim at presenting a detailed treatment of these cases.

384 4. Numerical studies

385 To illustrate the theoretical results of Section 3, we consider two toy examples: the first tar-
 386 get distribution considered is a normal-normal hierarchical model in which the components
 387 X_2, \dots, X_n are related through their mean, while the second one is a gamma-normal hierar-
 388 chical model in which X_2, \dots, X_n are related through their variance. In both cases, we show
 389 how to compute the optimal variance $\hat{\ell}$. We also study the performance of RWM samplers
 390 and conclude that even in relatively low-dimensional settings, the samplers behave according
 391 to the asymptotic results previously detailed.

392 4.1. Normal-normal hierarchical distribution

393 Consider an n -dimensional hierarchical target such that $X_1 \sim \mathcal{N}(0, 1)$ and $X_i|X_1 \sim \mathcal{N}(X_1, 1)$
 394 for $i = 2, \dots, n$. To sample from this distribution, we use a RWM algorithm with a
 395 $\mathcal{N}(\mathbf{x}, \ell^2 I_n/n)$ proposal distribution. This simple target shall relate Theorem 2 to the theo-
 396 retical results derived in [3].

397 Standard calculations lead to $X_1|\mathbf{X}_{2:n} \sim \mathcal{N}(\sum_{i=2}^n X_i/n, 1/n)$; as $n \rightarrow \infty$, $\mathbb{V}(X_1|\mathbf{X}_{2:n}) \rightarrow 0$
 398 almost surely. If we let $\mu_n = \sum_{i=2}^n X_i/n$ and $\tilde{X}_1 = n^{1/2}(X_1 - \mu_n)$, then $\tilde{X}_1|\mathbf{X}_{2:n} \sim \mathcal{N}(0, 1)$.
 399 Furthermore, the term $\sum_{i=2}^n (X_i - \mu_n - \tilde{X}_1/\sqrt{n})^2/n$ is reexpressed as $\sum_{i=2}^n (X_i - X_1)^2/n =$
 400 $\sum_{i=2}^n Z_i^2/n$ and thus converges in probability to $\mathbb{E}[Z^2] = \int (\frac{\partial}{\partial x} \log f(x|\mu))^2 f(x|\mu) dx = 1$,
 401 where Z_1, \dots, Z_n denote independent standard normal random variables. By Theorem 1, we
 402 can thus affirm that the component \tilde{X}_1 asymptotically behaves according to a one-dimensional
 403 RWM algorithm with a standard normal target and acceptance function as in (4); these do
 404 not, in the current case, depend on \mathbf{x} .

405 Evaluating the function $\gamma(x_1, z_1)$ in (7) is a simple task and leads to $\gamma(x_1, z_1) = z_1^2 + 1$. The
 406 AOS value is then found by maximizing

$$v(\ell) = 2\ell^2 \mathbb{E}_{Z_1} \left[\Phi \left(-\frac{\ell}{2} \sqrt{Z_1^2 + 1} \right) \right]$$

407 with respect to ℓ , where $Z_1 \sim \mathcal{N}(0, 1)$. This yields an AOS of $\hat{\ell}^2 = 4.00$ and a corresponding
 408 AOAR of $v(\hat{\ell})/\hat{\ell}^2 = 0.205$. These values are naturally smaller than those obtained for a
 409 target with i.i.d. components (5.66 and 0.234, respectively); indeed, the proposal distribution
 410 is formed of i.i.d. components and accordingly better suited for similar targets. Relying on a
 411 proposal with correlated components would however require a certain understanding of the
 412 target correlation structure, which goes against the general framework we wish to consider.

413 It is worth pointing out that the speed measure of the limiting diffusion process does not
 414 depend on X_1 in the present case. This holds for arbitrary densities f_1 and f satisfying the
 415 conditions in Section 2.1, provided that X_1 is a location parameter for X_i ($i \geq 2$). Since a
 416 variation in the location parameter does not perturb the roughness of the distribution, the
 417 AOS and AOAR found are valid both locally and globally. This means that $\hat{\ell}$, which remains
 418 fixed across iterations, is the best possible proposal scaling conditionally on the last position
 419 of the component X_1 (*i.e.* $\hat{\ell} = \hat{\ell}(x_1)$).

420 A second peculiarity of this example is that the target distribution is jointly normal with
 421 mean $\mathbf{0}$ and $n \times n$ covariance matrix Σ_n given by $\sigma_1^2 = 1$, $\sigma_j^2 = 2$ ($j = 2, \dots, n$), and
 422 $\sigma_{i,j} = 1 \forall i \neq j$ ($i, j = 1, \dots, n$). Normal distributions being invariant under orthogonal
 423 transformations, we can find a transformation under which the target components become
 424 mutually independent. The covariance matrix Σ_n is thus transformed into a diagonal matrix
 425 whose diagonal elements consist in the eigenvalues of Σ_n . In moderate to large dimensions,
 426 the eigenvalues can be approximated by $1/(n+1), (n+1), 1, \dots, 1$. It turns out that the
 427 optimal scaling problem for target distributions of this sort (*i.e.* formed of components that
 428 are i.i.d. up to a scaling term) has been studied in [1]. Solving for the AOS value and AOAR
 429 of the transformed target using Theorem 1 and Corollary 2 in [3] leads to values that are
 430 consistent with those obtained using Theorem 2 in Section 3.

431 To illustrate these theoretical results, we consider the 20-dimensional normal-normal target
 432 described above and run 50 RWM algorithms that differ by their proposal variance only.

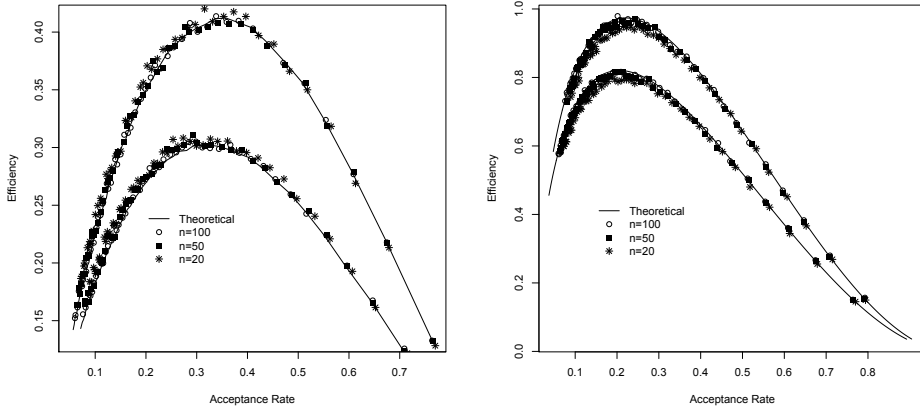


Figure 1: Efficiency of RWM algorithm against acceptance rate for the normal-normal hierarchical target. Left: efficiency of X_1 only; the top set of curves corresponds to homogeneous proposal variances. Right: efficiency of all n components; the top set of curves now corresponds to inhomogeneous proposal variances.

433 For each sampler, we perform 100,000 iterations (sufficient for convergence according to the
 434 autocorrelation function) and measure efficiency by recording the average squared jumping
 435 distance

$$\text{ASJD} = \frac{1}{N} \sum_{j=1}^N \sum_{i=1}^n \left(x_i^{(n)}[j] - x_i^{(n)}[j-1] \right)^2; \quad (8)$$

436 here, N is the number of iterations and n is the dimension of the target distribution. We
 437 also record the average acceptance rate of each algorithm, expressed as

$$\text{AAR} = \frac{1}{N} \sum_{j=1}^N \mathbb{1}\{\mathbf{x}^{(n)}[j] \neq \mathbf{x}^{(n)}[j-1]\}.$$

438 We repeat these steps for 50- and 100-dimensional normal-normal targets, and combine all
 439 three curves of efficiency versus acceptance rate on a graph along with the theoretical ef-
 440 ficiency curve of $v(\ell)$ versus the expected acceptance rate $v(\ell)/\ell^2$ (Figure 1, right graph,
 441 bottom set of curves). To assess the limiting behaviour of the coordinate X_1 , we also plot
 442 the ASJD of this single component (for the 20-, 50-, and 100-dimensional cases) along with
 443 the ASJD for the limiting one-dimensional RWM sampler described in Theorem 1 (Figure 1,
 444 left graph, top set of curves).

445 We now repeat the numerical experiment by taking advantage of the available target vari-
 446 ances in the tuning of the proposal distribution. Specifically, we let $Y_1 \sim \mathcal{N}(x_1, \ell^2/2n)$ be
 447 independent of $\mathbf{Y}_{2:n} \sim \mathcal{N}(\mathbf{x}_{2:n}, \ell^2/n)$ and run the RWM algorithm in dimensions 20, 50, and
 448 100. The resulting simulated and theoretical efficiency curves are illustrated in Figure 1 (left
 449 graph, bottom set of curves; right graph, top set of curves). Although efficiency curves for
 450 X_1 are lower when using inhomogeneous proposal variances, this approach still results in a
 451 better overall performance (the curves in the right graph are higher than with homogeneous
 452 variances). The optimized theoretical efficiency is 0.974, which is related to an AOAR of

453 0.221. Despite the fact that Theorems 1 and 2 are valid asymptotically, the simulation study
 454 yields efficiency curves that are very close together; the theorems thus seem applicable in
 455 relatively low-dimensional settings.

456 Each set of curves on the right graph of Figure 1 agrees about the optimal acceptance rates
 457 0.205 and 0.221, respectively. These optimal rates have been obtained by running an ho-
 458 mogeneous sampler with optimal variance $\hat{\ell}^2/n = 4/n$ and an inhomogeneous sampler with
 459 optimal variance $4.4/n$, each optimizing (6). Any other proposal variance leads to a point
 460 that is lower on the efficiency curve.

461 According to the shape of these curves, tuning the acceptance rate anywhere between 0.15
 462 and 0.3 would yield a loss of at most 10% in efficiency, and would still result in a Markov chain
 463 that rapidly explores its state space; in particular, using the usual 0.234 for this target would
 464 yield an almost optimal algorithm. Beyond finding the exact AOAR for a specific target
 465 distribution, there is thus a need for understanding when and why AOARs significantly differ
 466 from 0.234. At the present time, the only way to answer this question is by solving the
 467 optimal scaling problem for target distributions of interest.

468 4.2. Gamma-normal hierarchical distribution

469 As a second example, consider a gamma-normal hierarchical target such that $X_1 \sim \Gamma(\alpha, \lambda)$
 470 and $X_i|X_1 \sim \mathcal{N}(0, 1/X_1)$, $i = 2, \dots, n$. Although X_i ($i \geq 2$) are still normally distributed,
 471 the coordinate X_1 now acts through the variance of the normal variables. This results in
 472 a target that significantly differs from the distribution considered in the previous section,
 473 falling slightly outside the framework of Section 2 ($\frac{\partial}{\partial x_1} \log f(x|x_1)$ is now only locally Lipschitz
 474 continuous). We run the usual RWM algorithm to obtain a sample from this distribution.

475 Standard calculations lead to $X_1|\mathbf{X}_{2:n} \sim \Gamma(\alpha + (n-1)/2, \lambda + \sum_{i=2}^n X_i^2/2)$ and as $n \rightarrow$
 476 ∞ , $\mathbb{V}(X_1|\mathbf{X}_{2:n}) \rightarrow_p 0$. The WLLN-type expression in Theorem 1 may be reexpressed as
 477 $\sum_{i=2}^n (\mu_n + \tilde{X}_1/\sqrt{n})^2 X_i^2/n = (\mu_n + \tilde{X}_1/\sqrt{n})(\sum_{i=2}^n Z_i^2/n)$, where $\mathbf{Z}_{1:n}$ are independent standard
 478 normal random variables. The condition is thus satisfied as it converges in probability to
 479 $\mu(\underline{\mathbf{X}}) = \mathbb{E}_X[(\frac{\partial}{\partial \underline{\mathbf{X}}} \log f(X|\mu))^2]$. Using Stirling's formula, it is not difficult to show that the
 480 density of $\tilde{X}_1|\mathbf{X}_{2:n}$ converges almost surely to that of a $\mathcal{N}(0, 2/\mu^2(\underline{\mathbf{X}}))$. By Theorem 1,
 481 the coordinate \tilde{X}_1 asymptotically behaves according to an atypical one-dimensional RWM
 482 algorithm with a normal target; the target variance however varies from one iteration to the
 483 next, and so does the acceptance function in (4).

484 To optimize the efficiency of the algorithm, we analyze the speed measure in (6); in the
 485 present case, it is expressed as

$$v(\ell, x_1) = 2\ell^2 \mathbb{E}_{Z_1} \left[\Phi \left(-\frac{\ell}{2} \sqrt{\frac{1}{2} \frac{Z_1^2}{x_1^2} + x_1} \right) \right],$$

486 where $Z_1 \sim \mathcal{N}(0, 1)$. Maximizing the function $\mathbb{E}_{X_1}[v(\ell, X_1)]$ in Corollary 6 with respect to
 487 ℓ leads to the global AOS value, which is fixed across iterations; when $(\alpha, \lambda) = (3, 1)$ for
 488 instance, we find $\hat{\ell}^2 = 2.40$ and $\text{AOAR} = 0.204$.

489 The simulation study described in Section 4.1 has been performed for the gamma-normal
 490 target model with various α and λ . Specifically, for fixed α, λ , we consider a 10-dimensional

Table 1: Optimal efficiency and acceptance rate of chains in various dimension ($n = 10, 20, 50$), for different parameters α, λ of the gamma distribution for X_1 . The theoretical optimal efficiency and acceptance rate are also included for comparison.

Parameters	Optimal efficiency				Optimal acceptance rate			
	Theoretical	$n = 10$	$n = 20$	$n = 50$	Theoretical	$n = 10$	$n = 20$	$n = 50$
$\alpha = 2, \lambda = 1$	0.6381	0.6036	0.6246	0.6456	0.1934	0.1984	0.1968	0.1857
$\alpha = 2, \lambda = 2$	0.8169	0.7430	0.7862	0.8239	0.1815	0.1888	0.1759	0.1886
$\alpha = 2, \lambda = 3$	0.8420	0.7623	0.8170	0.8608	0.1517	0.1682	0.1527	0.1593
$\alpha = 3, \lambda = 1$	0.4889	0.4503	0.4736	0.4926	0.2037	0.2370	0.2158	0.2001
$\alpha = 3, \lambda = 2$	0.7541	0.6739	0.7139	0.7405	0.2038	0.2265	0.2233	0.2040
$\alpha = 3, \lambda = 3$	0.8648	0.7554	0.8075	0.8497	0.1922	0.1931	0.1930	0.1882

491 gamma-normal target distribution and run 50 RWM algorithms possessing their own proposal
 492 variance. For each sampler, we perform 1,000,000 iterations (again sufficient for convergence
 493 according to the autocorrelation function) and measure efficiency by recording the ASJD of
 494 each chain. We then repeat these steps for 20- and 50-dimensional targets. Table 1 presents
 495 the optimal efficiency and acceptance rate for various α, λ . Those results are compared to
 496 the theoretical optimal values obtained by maximizing $\mathbb{E}_{X_1}[v(\ell, X_1)]$.

497 Although the corresponding graphs are omitted here, they yield curves similar to those ob-
 498 tained in Figure 1 for the normal-normal target. We note that even if the gamma-normal
 499 departs from a jointly normal distribution assumption and does not yield as nice a target
 500 distribution as in the previous example, the AOAR obtained is not too far from the 0.234
 501 found for i.i.d. targets. The AOAR however tends to decrease as λ increases (*e.g.* 0.152 for
 502 $(\alpha, \lambda) = (2, 3)$).

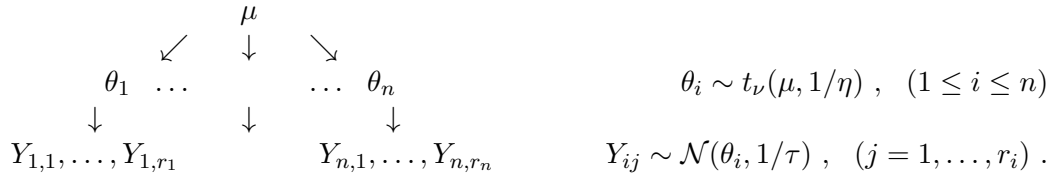
503 In the current example, it also turns out that the agreement between theoretical and sim-
 504 ulation results is altered for some values (α, λ) . As mentioned above, one of the Lipschitz
 505 conditions is only valid locally and so the change in $\frac{\partial}{\partial x_1} \log f(x|x_1)$ becomes arbitrarily steep
 506 as $X_1 \rightarrow 0$. The amplitude of X_1 movements is, therefore, not adequately controlled for some
 507 choices of (α, λ) that yield a density f_1 assigning a significant probability close to 0. In cases
 508 where regularity assumptions are not all satisfied, the applicability of theoretical results may
 509 thus be affected by the choice of hyperparameters.

510 5. Applications in Bayesian contexts

511 The theoretical results presented in this paper have wide applicability and may be used
 512 to improve not only RWM algorithms, but other samplers as well (RWM-within-Gibbs, for
 513 instance). The examples below study the performance of optimally tuned samplers in the
 514 context of hierarchical Bayesian models. They show that the RWM-within-Gibbs sampler
 515 with local variances (*i.e.* variances that are a function of the current state of the chain) is
 516 superior to its counterpart with a fixed variance. It is also superior to traditional RWM
 517 algorithms and even Adaptive Metropolis (AM) samplers, which use the history of the chain
 518 to recursively update the covariance matrix of their proposal distribution (see [11]).

519 *5.1. Scottish secondary school scores*

520 The dataset `ScotsSec` in the package `mlmRev` in R contains the scores attained by 3,435
 521 Scottish secondary school students on a standardized test taken at age 16. The primary
 522 schools attended by students are also recorded in this dataset; there are $n = 148$ different
 523 primary schools, and the number of students per primary school varies between 1 and 72.
 524 We use the following multilevel Bayesian framework to model these data



525 In this model, the variables $\mathbf{y}_{i,1:r_i}$ represent the observed scores obtained by the r_i students
 526 having attended primary school i , $i = 1, \dots, 148$. These observations are modeled according
 527 to a normal distribution with mean θ_i and variance $1/\tau$. The group sizes range from $r_{148} = 1$
 528 to $r_{61} = 72$. The variables $\boldsymbol{\theta}_{1:148}$, which represent the mean scores of the standardized test
 529 for students having attended each of the 148 primary schools, are modeled using a Student
 530 distribution with $\nu = 4$ degrees of freedom. A translated and scaled Student distribution
 531 $t_\nu(\mu, 1/\eta)$ has a density proportional to $[1 + \eta(x - \mu)^2/\nu]^{-(\nu+1)/2}$. The mean and precision
 532 of the Student distribution, along with the precision of the normally distributed data, are
 533 attributed non-informative priors: $\pi(\mu) \propto 1$, $\pi(\eta) \propto \eta^{-1}$, and $\pi(\tau) \propto \tau^{-1}$.

534 This model leads to the $(n + 3)$ -dimensional posterior density

$$\begin{aligned}
 \pi(\mu, \eta, \tau, \boldsymbol{\theta}_{1:n} | \{Y_{ij}\}) &\propto \eta^{-1} \tau^{-1} \prod_{i=1}^n \sqrt{\eta} \left[1 + \frac{\eta(\theta_i - \mu)^2}{\nu} \right]^{-(\nu+1)/2} \\
 &\quad \prod_{i=1}^n \prod_{j=1}^{r_i} \sqrt{\tau} \exp \left\{ -\frac{\tau}{2} (y_{ij} - \theta_i)^2 \right\}. \quad (9)
 \end{aligned}$$

535 The posterior density is too complex for analytic computation, and numerical integration
 536 must be ruled out due to the dimensionality of the problem. This distribution is best sampled
 537 with MCMC methods, although a classical Gibbs sampler must be ruled out, as the Student
 538 distribution destroys conjugacy. In the current setting, we propose to use a RWM-within-
 539 Gibbs with four blocks of variables: μ , η , τ , and $\boldsymbol{\theta}_{1:n}$. We are also interested in assessing
 540 the performance of full-dimensional RWM and AM algorithms in which μ , η , τ , and $\boldsymbol{\theta}_{1:n}$ are
 541 updated at once.

542 The RWM-within-Gibbs performs one-dimensional updates of μ , η , and τ using target densi-
 543 ties $f(\mu | \eta, \tau, \boldsymbol{\theta}_{1:n}, \{Y_{ij}\})$, $f(\eta | \mu, \tau, \boldsymbol{\theta}_{1:n}, \{Y_{ij}\})$, and $f(\tau | \mu, \eta, \boldsymbol{\theta}_{1:n}, \{Y_{ij}\})$. It then performs an
 544 n -dimensional update of $\boldsymbol{\theta}_{1:n}$ with respect to the conditional density $f(\boldsymbol{\theta}_{1:n} | \mu, \eta, \tau, \{Y_{ij}\}) =$
 545 $\prod_{i=1}^n f(\theta_i | \mu, \eta, \tau, \mathbf{Y}_{i,1:r_i})$.

546 Since each block of variables is updated individually using a RWM sampler, we may compute
 547 local proposal variances for the fourth block using (6) and (7) in Theorem 2. The proposal
 548 variances maximizing (6) are adjusted according to the roughness of their corresponding

549 target component's distribution, and should offer a better performance than a fixed proposal
 550 variance.

551 The target distribution of the fourth block satisfies

$$f(\boldsymbol{\theta}_{1:n}|\mu, \eta, \tau, \{Y_{ij}\}) \propto \prod_{i=1}^n \left[1 + \frac{\eta(\theta_i - \mu)^2}{\nu} \right]^{-\frac{\nu+1}{2}} \exp \left\{ -\frac{\tau}{2} \sum_{j=1}^{r_i} (y_{ij} - \theta_i)^2 \right\},$$

552 hence the partial derivative of the one-dimensional log density with respect to θ_i is

$$\frac{\partial}{\partial \theta_i} \log f(\theta_i|\mu, \eta, \tau, \mathbf{Y}_{i,1:r_i}) = \tau \sum_{j=1}^{r_i} (y_{ij} - \theta_i) - \frac{\nu+1}{\nu} \sqrt{\eta} \left(\frac{T_i}{1 + T_i^2/\nu} \right), \quad (10)$$

553 where $T_i = \sqrt{\eta}(\theta_i - \mu) \sim t_\nu(0, 1)$, $i = 1, \dots, n$. Since the variables μ, η, τ are updated
 554 separately, then the first term in (7) is null, leading to

$$\gamma_i(\mu, \eta, \tau) = \mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} f(\theta_i|\mu, \eta, \tau, \mathbf{Y}_{i,1:r_i}) \right)^2 \right]. \quad (11)$$

555 Optimizing (6) leads to local, inhomogeneous proposal variances of the form $2.38^2/\{n\gamma_i(\mu, \eta, \tau)\}$.

556 The terms $\gamma_i(\mu, \eta, \tau)$ in the proposal variances are not easy to obtain explicitly as the ex-
 557 pectation in (11) must be computed with respect to the conditional distribution of θ_i given
 558 $(\mu, \eta, \tau, \mathbf{Y}_{i,1:r_i})$, which is not a Student distribution anymore. However, the terms $\gamma_i(\mu, \eta, \tau)$
 559 may be averaged over the random variables $\mathbf{Y}_{i,1:r_i}$. Squaring (10) and computing the expect-
 560 ation first with respect to $\mathbf{Y}_{i,1:r_i}$ and then with respect to θ_i easily leads to

$$\mathbb{E}[\gamma_i(\mu, \eta, \tau)] = r_i \tau + \eta \frac{(\nu+1)^2}{\nu(\nu+2)} \frac{\Gamma((\nu+1)/2) \Gamma((\nu+4)/2)}{\Gamma(\nu/2) \Gamma((\nu+5)/2)}.$$

561 These terms yield local proposal variances that have been averaged over all possible datasets;
 562 these are the best local variances for the model under study when no information about the
 563 observations is available.

564 The RWM-within-Gibbs is then implemented using Gaussian proposal distributions with $\sigma_1 =$
 565 0.95 , $\sigma_2 = 0.025$, and $\sigma_3 = 0.0005$ for μ , η , and τ . This yields acceptance rates in the range
 566 35%-50% for each sub-algorithm, as prescribed in the literature for one-dimensional target
 567 distributions (see [18]). We update $\boldsymbol{\theta}_{1:148}$ using a Gaussian proposal with local variances
 568 $2.38^2/\{n\mathbb{E}[\gamma_i(\mu, \eta, \tau)]\}$.

569 These steps are then repeated by running a RWM-within-Gibbs in which $\boldsymbol{\theta}_{1:148}$ is updated
 570 using a fixed proposal variance of 5^2 . We also run a 151-dimensional RWM sampler with
 571 a $\mathcal{N}((\mu, \eta, \tau, \boldsymbol{\theta}_{1:148}), 4^2/151 * \text{diag}(1, 0.01, 0.001, 1, \dots, 1))$ proposal distribution, and an AM
 572 algorithm in which the tuning factor of the proposal covariance matrix is 8.

573 The ASJD of the chain in (8) offers a reliable way of comparing the four samplers; it is reported
 574 in the first column of Table 2. A large value of this measure (relative to other samplers) is
 575 indicative of a process that rapidly explores its space, and is equivalent to ordering samplers
 576 according to their lag-1 autocorrelations. We also compare the relative efficiency of these
 577 samplers by calculating the effective sample size (ESS) of the variables μ , η , τ , and θ_2 . The

Table 2: *Scottish dataset*: Efficiency and time-adjusted efficiency measures for the four samplers tested.

Sampler	Efficiency			Time-adjusted efficiency	
	Mean ASJD (a)	Min ESS (e)	Mean time (s)	a/s (×100)	e/s (×100)
RWM	2.9712	74.30	145.50	2.0421	51.07
Fixed RWM-w-G	6.4279	157.09	147.77	4.3499	106.31
Local RWM-w-G	8.4108	272.70	148.06	5.6807	184.18
Adapt. Met.	5.2476	473.83	1,081.52	0.4852	43.81

578 effective sample size represents the number of uncorrelated samples that are produced from
579 the output of the sampler. It is also used as a convergence diagnostic: when its value is too
580 small (< 100), we may have reasonable doubts that the chain really has converged. It is
581 computed as

$$ESS = \frac{N}{1 + 2 \sum_{k=1}^{\infty} \gamma(k)},$$

582 where N is the number of samples and $\sum_{k=1}^{\infty} \gamma(k)$ is the sum of lag- k sample autocorrelations.
583 An ESS is produced for each variable; since we want to measure the number of samples that
584 are uncorrelated over all variables, we report the minimum ESS (2nd column of Table 2).
585 The ASJD and minimum ESS values are averaged over 10 runs of 100,000 iterations each,
586 with a burn-in period of 1,000. These quantities are then normalized relative to the average
587 running time of samplers (3rd column); this respectively yields the average square jumping
588 distance per second (4th column), and the number of uncorrelated samples generated every
589 second (5th column).

590 According to these results, the RWM-within-Gibbs with local variances is 1.3 times more
591 efficient than the one with a fixed variance; the efficiency gain is even greater (1.7) if we
592 consider the minimum ESS instead of the ASJD. Although the RWM sampler offers a slight
593 improvement in terms of running time, it still results in efficiency measures that are sig-
594 nificantly smaller than those of the RWM-within-Gibbs. The Adaptive Metropolis sampler
595 could be an interesting alternative to the RWM-within-Gibbs, if it were not as expensive
596 in terms of computational resources. Indeed, even if its ASJD is smaller than that of the
597 RWM-within-Gibbs, its minimum ESS is greater. This sampler however requires significantly
598 more time than the other samplers to complete its 100,000 iterations. When correcting for
599 computational effort, it thus badly loses ground to its competitors.

600 The results in Table 2 thus illustrate that there is an important efficiency gain that is available
601 from preferring a local RWM-within-Gibbs over its constant counterpart. Given that running
602 times for both approaches are equivalent, we should clearly use local proposal variances
603 whenever possible.

604 5.2. Stochastic volatility model

605 As a second example, we wish to study the performance of MCMC samplers in the context
606 of a Bayesian hierarchical model that does not respect the regularity assumptions imposed

607 by the theory of Section 3. We consider a stochastic volatility model in which the latent
 608 volatilities form an order-1 autoregressive process. The model, similar to those studied in
 609 [10] and [13], expresses the mean corrected returns d_i and log volatilities X_i , for $i \geq 1$, as

$$\begin{aligned} d_i &= \varepsilon_i \exp\{X_i/2\}, \\ X_{i+1} &= \phi X_i + \eta_{i+1}. \end{aligned}$$

610 The variables $\varepsilon_i \sim \mathcal{N}(0, 1)$ and $\eta_i \sim \mathcal{N}(0, \tau^2)$ are uncorrelated white noises and we set $X_1 \sim$
 611 $\mathcal{N}(0, \tau^2/(1-\phi^2))$. Priors for the parameters τ^2 and ϕ are $\tau^2 \sim \text{IG}(\delta, \lambda)$ and $(\phi+1)/2 \sim \beta(a, b)$,
 612 where $\text{IG}(\delta, \lambda)$ is the inverse gamma distribution with density proportional to $x^{-(\delta+1)}e^{-\lambda/x}$.
 613 This model leads to an $(n+2)$ -dimensional posterior density $\pi(\tau^2, \phi, X_1, \dots, X_n | \mathbf{d}_{1:n})$.

614 Before pursuing the analysis, we note that τ^2 and ϕ are constrained to subsets of \mathbb{R} ; since the
 615 target density is rather sensitive to changes in these parameters, this will potentially affect
 616 the performance of MCMC approaches. To ensure fluidity in the samplers implemented, we
 617 apply the transformations $\tau^2 = \exp\{\kappa\}$ and $\phi = \tanh(\omega)$. The new variables κ, ω take values
 618 in \mathbb{R} and the resulting $(n+2)$ -dimensional posterior density is given by

$$\begin{aligned} \pi(\kappa, \omega, \mathbf{x}_{1:n} | \mathbf{d}_{1:n}) &\propto \exp\left\{-\kappa\left(\frac{n}{2} + \delta\right)\right\} \frac{e^{-\omega(2b+1)}}{(1+e^{-2\omega})^{a+b+1}} \exp\left\{-\frac{1}{2} \sum_{i=1}^n (x_i + d_i^2 e^{-x_i})\right\} \\ &\times \exp\left\{-\frac{e^{-\kappa}}{2} \left[2\lambda + \frac{4e^{-2\omega}}{(1+e^{-2\omega})^2} x_1^2 + \sum_{i=2}^n \left(x_i - \left(\frac{1-e^{-2\omega}}{1+e^{-2\omega}}\right)x_{i-1}\right)^2\right]\right\}. \end{aligned}$$

619 Using a 100-dimensional dataset $\mathbf{d}_{1:100}$ exhibiting low correlation (obtained from the stochas-
 620 tic volatility model with $\phi = 0.1$ and $\tau^2 = 0.75$), we sample this posterior density using
 621 RWM-within-Gibbs (local and fixed variances), traditional RWM, and AM algorithms. Hy-
 622 perparameters are set to $\delta = 1$, $\lambda = 0.75$, $a = 10$, and $b = 6$.

623 For the RWM-within-Gibbs, we propose to divide the variables into 3 blocks: κ , ω , and $\mathbf{X}_{1:n}$.
 624 The proposal standard deviations associated to κ and ω are set to 0.2 and 0.27 respectively;
 625 each sub-algorithm thus accepts candidates according to a rate of $\approx 45\%$. The n -dimensional
 626 update of $\mathbf{X}_{1:n}$ is performed according to the conditional target density $\pi(\mathbf{x}_{1:n} | \kappa, \omega, \mathbf{d}_{1:n})$. In
 627 the case of the RWM-within-Gibbs with local variances, the terms

$$\gamma_i(\kappa, \omega) = \mathbb{E}\left[\left(\frac{\partial}{\partial X_i} \log \pi(\mathbf{X}_{1:n} | \kappa, \omega, \mathbf{d}_{1:n})\right)^2\right], \quad i = 1, \dots, n$$

628 in (7) are not easy to obtain as the full conditional distribution (given the data) is not nor-
 629 mally distributed anymore. As before, we solve this problem by computing the expectation
 630 above with respect to $\mathbf{d}_{1:n}$ first, and then with respect to $\mathbf{X}_{1:n}$. The resulting proposal vari-
 631 ances are thus averaged over all possible datasets; they are the best local proposal variances,
 632 independently of the specific dataset considered. Optimizing (6) for $i = 1, \dots, n$ yields the
 633 n -dimensional vector

$$\frac{2.38^2}{n} \left(\frac{1}{2} + e^{-\kappa}, \frac{1}{2} + e^{-\kappa} \left(1 + \left(\frac{1-e^{-2\omega}}{1+e^{-2\omega}}\right)^2\right), \dots, \frac{1}{2} + e^{-\kappa} \left(1 + \left(\frac{1-e^{-2\omega}}{1+e^{-2\omega}}\right)^2\right), \frac{1}{2} + e^{-\kappa}\right)^{-1}. \quad (12)$$

634 For the RWM-within-Gibbs with a fixed proposal variance, the proposal standard deviations
 635 associated to κ and ω are still 0.2 and 0.27. We then use the theory of Section 3 to obtain

Table 3: *Stochastic volatility* - Efficiency and time-adjusted efficiency measures for the four samplers tested.

Sampler	Efficiency			Time-adjusted efficiency	
	Mean ASJD (<i>a</i>)	Min ESS (<i>e</i>)	Mean time (<i>s</i>)	<i>a/s</i> (× 1,000)	<i>e/s</i> (× 1,000)
RWM	0.3994	103.37	367.34	1.0873	281.40
Fixed RWM-w-G	0.6420	116.58	371.93	1.7261	313.45
Local RWM-w-G	0.6740	132.40	371.38	1.8149	356.51
Adapt. Met.	0.6320	347.76	1,149.26	0.5499	302.59

636 an approximately optimal acceptance rate of 0.2 for the block $\mathbf{X}_{1:n}$. We reach a similar con-
 637 clusion for the traditional RWM sampler. Naturally, we have to keep in mind that regularity
 638 assumptions are violated in the current context; the theoretical results might not be robust
 639 to a departure from those assumptions. In fact, given that the X_i s are correlated, we expect
 640 the Adaptive Metropolis sampler to better capture this design and to outdo its competitors.

641 The initial covariance matrix of the Adaptive Metropolis algorithm is the $(n + 2)$ -dimensional
 642 identity matrix. We tune its acceptance rate as close as possible to 0.234, as suggested in the
 643 literature. For each sampler, we average the ASJD and minimum ESS over 10 runs of 200,000
 644 iterations each, from which the first 10,000 iterations are discarded as burn-in. Time-adjusted
 645 ESJD and minimum ESS are again used as measures of efficiency; their values are reported
 646 in Table 3.

647 In terms of ASJD, the RWM-within-Gibbs with local variances is the best option, although its
 648 competitors also offer decent performances. The AM sampler does better, in absolute, for the
 649 minimum ESS; when accounting for computational effort however, the AM ends up outdone
 650 by the RWM-within-Gibbs (local and fixed). As before, we notice a net efficiency gain when
 651 preferring local variances to a fixed one in the RWM-within-Gibbs (net gain between 5%
 652 and 13%, depending on the efficiency measure). This modest gain is explained by the fact
 653 that, for the specific model studied, variations in κ and ω do not have a huge impact on
 654 the value of the local variances in (12). In spite of this, the impact of using local variances
 655 remains positive; generally, there does not seem to be a risk associated to using such local
 656 variances. Furthermore, the theoretical results seem applicable to contexts where regularity
 657 assumptions are violated (to some extent).

658 6. Discussion

659 In this paper, we have studied the tuning of RWM algorithms applied to single-level hierar-
 660 chical target distributions. The optimal variance of the Gaussian proposal distribution has
 661 been found to depend on a measure of roughness of the density f with respect to x as before,
 662 but also with respect to the mixing coordinate x_1 . This leads to local proposal variances that
 663 are a function of the mixing parameter x_1 . It is however possible to average over the random
 664 variable X_1 to find a globally optimal proposal variance. In the case where X_1 is a location
 665 parameter, it does not affect the roughness of the density f and the optimal proposal scaling

666 is valid both locally and globally.

667 Higher-level hierarchies could be studied using a similar approach. A target featuring p
668 mixing components, expressed as

$$\pi(\mathbf{x}) = \prod_{j=1}^p f_j(x_j) \prod_{i=p+1}^n f(x_i | \mathbf{x}_{1:p}) ,$$

669 with $\mathbf{x}_{1:p} = (x_1, \dots, x_p)$ would lead to a result similar to Theorem 2, but with the function

$$\gamma(\mathbf{x}_{1:p}, \mathbf{z}_{1:p}) = \mathbb{E}_X \left[\left(\sum_{j=1}^p z_j \frac{\partial}{\partial x_j} \log f(X | \mathbf{x}_{1:p}) \right)^2 \right] + \mathbb{E}_X \left[\left(\frac{\partial}{\partial X} \log f(X | \mathbf{x}_{1:p}) \right)^2 \right] ,$$

670 where $\mathbf{z}_{1:p} = (z_1, \dots, z_p)$ come from independent $\mathcal{N}(0, 1)$ random variables. For a tar-
671 get whose mixing component (X_p say) depends itself on higher-level mixing components
672 X_1, \dots, X_{p-1} , expressed as $\pi(\mathbf{x}) = f_1(\mathbf{x}_{1:p}) \prod_{i=p+1}^n f(x_i | x_p)$, the conclusions of Theorem
673 2 are still valid. These generalizations also hold for Corollary 7, with obvious adjustments
674 ($Z_1 \sim \mathcal{N}(0, \ell^2 \kappa_1^2/n), \dots, Z_p \sim \mathcal{N}(0, \ell^2 \kappa_p^2/n)$). Similar extensions may be derived for other
675 hierarchical models.

676 In the simulation study of Section 4, we found that the optimal acceptance rate most often
677 lies around 0.2. In the gamma-normal example, there were some values of α, λ that led to
678 significantly lower optimal acceptance rates (0.15 when $\alpha = 2, \lambda = 3$). The usual 0.234 is
679 thus quite robust and, if preferred, should lead to an efficient version of the sampler. In the
680 case of correlated targets, it would however be wiser to settle for an acceptance rate slightly
681 below 0.234. Since we investigate correlated targets with a proposal distribution featuring a
682 diagonal covariance matrix, it is not surprising to find an AOAR lower than 0.234; the latter
683 is the AOAR for exploring a target distribution with independent components, which is an
684 ideal situation when relying on a proposal distribution with independent components.

685 We conclude by outlining that the concept of locally optimal proposal variances reveals
686 itself to be of interest with other types of samplers, such as RWM-within-Gibbs algorithms.
687 Indeed, the asymptotic results of Section 3 are proof of the theoretical superiority of RWM-
688 within-Gibbs over RWM when sampling from hierarchical targets. The examples of Section 5
689 illustrate the efficiency gain from using a RWM-within-Gibbs with local variances over some
690 competitors, including an adaptive sampler. Similar ideas may also be applied to different
691 samplers such as Metropolis-adjusted Langevin algorithms (MALA), but this goes beyond
692 the scope of this paper.

693 **Conflicts of Interest**

694 The author declares that there is no conflict of interest regarding the publication of this
695 paper.

696 **Funding Statement**

697 This work was supported by the National Sciences and Engineering Research Council of
698 Canada (Grant number: RGPIN/03931-2014).

699 **Acknowledgements**

700 The author is grateful to P. Gagnon for useful comments. The author also wishes to thank
701 anonymous referees for very constructive comments that greatly improved the manuscript.

702 **References**

- 703 [1] M. Bédard, On the Robustness of Optimal Scaling for Random Walk Metropolis Algo-
704 rithms., Ph.D. thesis, University of Toronto, 2006.
- 705 [2] M. Bédard, Weak convergence of Metropolis algorithms for non-i.i.d. target distributions,
706 *Ann. Appl. Probab.* 17 (2007) 1222–1244.
- 707 [3] M. Bédard, Optimal acceptance rates for Metropolis algorithms: Moving beyond 0.234,
708 *Stoch. Process. Appl.* 118 (2008) 2198–2222.
- 709 [4] M. Bédard, Hierarchical models: Local proposal variances for RWM-within-Gibbs and
710 MALA-within-Gibbs, *Comput. Statist. Data Anal.* 109 (2017) 231–246.
- 711 [5] M. Bédard, R. Douc, E. Moulines, Scaling analysis of multiple-try MCMC methods,
712 *Stoch. Process. Appl.* 122 (2012) 758–786.
- 713 [6] A. Beskos, N. Pillai, G. Roberts, J. Sanz-Serna, A. Stuart, Optimal tuning of the hybrid
714 Monte Carlo algorithm, *Bernoulli* 19 (2013) 1501–1534.
- 715 [7] A. Beskos, G. Roberts, A. Stuart, Optimal scalings for local Metropolis-Hastings chains
716 on nonproduct targets in high dimensions, *Ann. Appl. Probab.* 19 (2009) 863–898.
- 717 [8] P. Billingsley, Convergence of probability measures, *Wiley Series in Probability and*
718 *Statistics: Probability and Statistics*, John Wiley & Sons Inc., New York, second edition,
719 1999.
- 720 [9] S. Ethier, T. Kurtz, *Markov processes: Characterization and Convergence*, *Wiley Series*
721 *in Probability and Mathematical Statistics: Probability and Mathematical Statistics*,
722 John Wiley & Sons Inc., New York, 1986.
- 723 [10] M. Girolami, B. Calderhead, Riemann manifold Langevin and Hamiltonian Monte Carlo
724 methods, *J. Roy. Statist. Soc. Ser. B* 73 (2011) 123–214.
- 725 [11] H. Haario, E. Saksman, J. Tamminen, An adaptive Metropolis algorithm, *Bernoulli* 7
726 (2001) 223–242.
- 727 [12] W. Hastings, Monte Carlo sampling methods using Markov chains and their applications,
728 *Biometrika* 57 (1970) 97–109.
- 729 [13] S. Kim, N. Shepard, S. Chib, Stochastic volatility: likelihood inference and comparison
730 with ARCH models, *The Review of Economic Studies* 65 (1998) 361–393.
- 731 [14] J. Mattingly, N. Pillai, A. Stuart, Diffusion limits of random walk Metropolis in high
732 dimensions, *Ann. Appl. Probab.* 22 (2012) 881–930.

- 733 [15] N. Metropolis, A. Rosenbluth, M. Rosenbluth, A. Teller, E. Teller, Equation of state
734 calculations by fast computing machines, *The Journal of Chemical Physics* 21 (1953)
735 1087–1092.
- 736 [16] P. Neal, G. Roberts, Optimal scaling for partially updating MCMC algorithms, *Ann.*
737 *Appl. Probab.* 16 (2006) 475–515.
- 738 [17] G. Roberts, A. Gelman, W. Gilks, Weak convergence and optimal scaling of random
739 walk Metropolis algorithms, *Ann. Appl. Probab.* 7 (1997) 110–120.
- 740 [18] G. Roberts, J. Rosenthal, Optimal scaling for various Metropolis-Hastings algorithms,
741 *Stat. Sci.* 16 (2001) 351–367.
- 742 [19] T. Xifara, C. Sherlock, S. Livingstone, S. Byrne, M. Girolami, Langevin diffusions and
743 the Metropolis-adjusted Langevin algorithm, *Stat. Probab. Lett.* 91 (2014) 14–19.

744 A. Appendix : Proofs of theorems

745 We now proceed to prove Theorems 1 and 2. To assess weak convergence of the processes
746 $\{\tilde{W}_1^{(n)}(t); t \geq 0\}$ and $\{W_2^{(n)}(t); t \geq 0\}$ in the Skorokhod topology, we first verify weak con-
747 vergence of finite-dimensional distributions. Whereas these processes are not themselves
748 Markovian, they are $\mathcal{F}^{\tilde{\mathbf{W}}^{(n)}}(t)$ -progressive and $\mathcal{F}^{(W_1^{(n)}, W_i^{(n)})}(t)$ -progressive \mathbb{R} -valued processes
749 respectively, and the aim of this section is to establish their convergence to some Markov pro-
750 cesses. According to Theorem 8.2 of Chapter 4 in [9], we thus look at the pseudo generator of
751 $\{\tilde{W}_1^{(n)}(t); t \geq 0\}$ (resp. $\{W_2^{(n)}(t); t \geq 0\}$), the univariate process associated to the component
752 X_1 (resp. X_2) in the rescaled RWM algorithm introduced at the end of Section 2. We then
753 verify \mathcal{L}^1 -convergence to the generator of the special RWM sampler with acceptance rule (4)
754 (resp. the generator of the diffusion in (5)).

755 To complete the proofs, Theorem 7.8 of Chapter 3 in [9] says that we must also assess the
756 relative compactness of $\{\tilde{W}_1^{(n)}(t); t \geq 0\}$ and $\{W_2^{(n)}(t); t \geq 0\}$ for $n = 2, 3, \dots$, as well as
757 the existence of a countable dense set on which the finite-dimensional distributions weakly
758 converge. This is achieved by using Corollary 8.6 of Chapter 4 in [9]; in the setting of Theorem
759 1, the satisfaction of applicability conditions is immediate; in the setting of Theorem 2, the
760 satisfaction of the first condition is immediate, while the verification of the second condition
761 is briefly discussed in Section A.2.

762 A.1. Proof of Theorem 1

763 In Theorem 1, it is assumed that $\{\tilde{W}_1^{(n)}(t); t \geq 0\}$ is the component of interest in $\{\tilde{\mathbf{W}}^{(n)}(t); t \geq$
764 $0\}$. Define the pseudo generator of $\{\tilde{W}_1^{(n)}(t); t \geq 0\}$ as

$$\tilde{G}_n h(\tilde{W}_1^{(n)}(t)) = \mathbb{E} \left[h(\tilde{W}_1^{(n)}(t+1)) - h(\tilde{W}_1^{(n)}(t)) \middle| \mathcal{F}^{\tilde{\mathbf{W}}^{(n)}}(t) \right],$$

765 where h is an arbitrary test function. By setting $\xi_n(t) = h(\tilde{W}_1^{(n)}(t))$ and $\varphi_n(t) = \tilde{G}_n h(\tilde{W}_1^{(n)}(t))$,
766 conditions in part (c) of Theorem 8.2 (Chap. 4 in [9]) reduce to $\mathbb{E} \left[\left| \tilde{G}_n h(\tilde{W}_1^{(n)}(t)) - \tilde{G} h(\tilde{W}_1(t)) \right| \right]$

767 $\rightarrow 0$ as $n \rightarrow \infty$ for $h \in \overline{C}$ (the space of continuous and bounded functions on \mathbb{R}), where
 768 $\tilde{G}h(\tilde{W}_1(t))$ is the generator of the special RWM sampler described in Theorem 1.

769 The above may be reexpressed as $\mathbb{E} \left[\left| \tilde{G}_n h(\tilde{X}_1) - \tilde{G}h(\tilde{X}_1) \right| \right] \rightarrow 0$ as $n \rightarrow \infty$, where

$$\tilde{G}_n h(\tilde{x}_1) = \mathbb{E}_{\tilde{Y}_1} \left[\left(h(\tilde{Y}_1) - h(\tilde{x}_1) \right) \mathbb{E}_{\mathbf{Y}_{2:n}} \left[\alpha(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right] \right]$$

770 with $\tilde{\mathbf{x}}^{(n)} = (\tilde{x}_1, x_2, \dots, x_n)$ and similarly for $\tilde{\mathbf{Y}}^{(n)}$. The density of $\tilde{\mathbf{x}}^{(n)}$ is $\frac{1}{\sqrt{n}} \pi(\mu_n + \frac{\tilde{x}_1}{\sqrt{n}}, \mathbf{x}_{2:n})$

771 with π as in (1), and thus $\alpha(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) = 1 \wedge \frac{\pi(\mu_n + \tilde{Y}_1/\sqrt{n}, \mathbf{Y}_{2:n})}{\pi(\mu_n + \tilde{x}_1/\sqrt{n}, \mathbf{x}_{2:n})}$; hereafter, $1 \wedge x = \min(1, x)$.

772 Furthermore,

$$\tilde{G}h(\tilde{x}_1) = \mathbb{E}_{\tilde{Y}_1} \left[\left(h(\tilde{Y}_1) - h(\tilde{x}_1) \right) \alpha^*(\tilde{x}_1, \tilde{Y}_1 | \mathbf{x}) \right]$$

773 with α^* as in (4). Note that there is a slight abuse of notation as, although h is a function
 774 of x_1 only, the generator $\tilde{G}_n h(\tilde{x}_1)$ is a function of $\tilde{\mathbf{x}}^{(n)}$; a similar remark holds for $\tilde{G}h(\tilde{x}_1)$.
 775 We now proceed to verify this condition. Hereafter, we use $\rightarrow_{a.s.}$, \rightarrow_p , and \rightarrow_d to denote
 776 convergence almost surely, in probability, and in distribution.

777 In the current context where there is no time-rescaling factor, the limiting process shall
 778 remain a RWM algorithm. For $h \in \overline{C}$ and some $K > 0$, the triangle inequality implies

$$\begin{aligned} \mathbb{E} \left[\left| \tilde{G}_n h(\tilde{X}_1) - \tilde{G}h(\tilde{X}_1) \right| \right] &\leq K \mathbb{E} \left[\left| \alpha(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) - \alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right| \right] & \text{(A.1)} \\ &+ K \mathbb{E} \left[\left| \alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) - \alpha_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right| \right] \\ &+ K \mathbb{E} \left[\left| \mathbb{E}_{\mathbf{Y}_{2:n}} \left[\alpha_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right] - \alpha^*(\tilde{X}_1, \tilde{Y}_1 | \mathbf{X}) \right| \right], \end{aligned}$$

779 where the function $\alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)})$ shall be defined in Lemma B.1 and $\alpha_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) =$
 780 $1 \wedge \exp \left\{ \varepsilon_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right\}$. Here,

$$\begin{aligned} \varepsilon_1(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) &= \log \frac{f_1(\mu_n + \frac{\tilde{Y}_1}{\sqrt{n}} | \mathbf{x}_{2:n})}{f_1(\mu_n + \frac{\tilde{x}_1}{\sqrt{n}} | \mathbf{x}_{2:n})} + \sum_{i=2}^n \frac{\partial}{\partial x} \log f(x | \mu_n + \frac{\tilde{x}_1}{\sqrt{n}}) \Big|_{x=x_i} (Y_i - x_i) \\ &\quad - \frac{\ell^2}{2n} \sum_{i=2}^n \left(\frac{\partial}{\partial x} \log f(x | \mu_n + \frac{\tilde{x}_1}{\sqrt{n}}) \Big|_{x=x_i} \right)^2, \end{aligned} \quad \text{(A.2)}$$

781 with $\frac{1}{\sqrt{n}} f_1(\mu_n + \frac{\tilde{x}_1}{\sqrt{n}} | \mathbf{x}_{2:n})$ representing the conditional density of \tilde{X}_1 given $\mathbf{x}_{2:n}$.

782 By Lemmas B.1 and B.2, the first and second terms in (A.1) respectively converge to 0 as
 783 $n \rightarrow \infty$; in the sequel, we thus study the last term. Since $\mathbf{Y}_{2:n} \sim \mathcal{N}(\mathbf{x}_{2:n}, \ell^2 I_{n-1}/n)$, the
 784 second and third terms on the right of (A.2) are normally distributed with mean M and
 785 variance V , where $V = -2M = \frac{\ell^2}{n} \sum_{i=2}^n \left(\frac{\partial}{\partial x} \log f(x | \mu_n + \frac{\tilde{x}_1}{\sqrt{n}}) \Big|_{x=x_i} \right)^2$.

786 By assumption, this variance term converges in probability to $\ell^2 \tilde{\gamma}(\mu)$; hence, the last two
 787 terms on the right of (A.2) converge in probability to a $\mathcal{N}(-\ell^2 \tilde{\gamma}(\mu)/2, \ell^2 \tilde{\gamma}(\mu))$. Regularity
 788 conditions allow us to invoke the (multivariate) Continuous Mapping Theorem, which implies

$$\alpha_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \rightarrow_p 1 \wedge \exp \left\{ \mathcal{N} \left(\log \frac{g_1(\tilde{Y}_1 | \mathbf{X})}{g_1(\tilde{X}_1 | \mathbf{X})} - \frac{\ell^2}{2} \tilde{\gamma}(\mu), \ell^2 \tilde{\gamma}(\mu) \right) \right\}.$$

789 Proposition 2.4 in [17] then claims that the expectation of $1 \wedge \exp\{Z\}$, where Z is the
790 normal random variable just introduced, is equal to $\alpha^*(\tilde{X}_1, \tilde{Y}_1 | \mathbf{X})$. The Bounded Convergence
791 Theorem can then be used to conclude that the last term in (A.1) converges to 0 as $n \rightarrow \infty$.

792 *A.2. Proof of Theorem 2*

793 In Theorem 2, it is assumed that $\{W_i^{(n)}(t); t \geq 0\}$ ($i = 2, \dots, n$) is the component of interest
794 in the rescaled process $\{\mathbf{W}^{(n)}(t); t \geq 0\}$. Without loss of generality, fix $i = 2$ and define the
795 pseudo generator of $\{W_2^{(n)}(t); t \geq 0\}$ as

$$G_n h(W_2^{(n)}(t)) = n \mathbb{E} \left[h(W_2^{(n)}(t + \frac{1}{n})) - h(W_2^{(n)}(t)) \mid \mathcal{F}^{(W_1^{(n)}, W_2^{(n)})}(t) \right],$$

796 where h is an arbitrary test function.

797 By setting $\xi_n(t) = h(W_2^{(n)}(t))$ and $\varphi_n(t) = G_n h(W_2^{(n)}(t))$, part (c) of Theorem 8.2 (Chapter
798 4 in [9]) reduces to the conditions $\sup_n \sup_{s \leq T} \mathbb{E}[|G_n h(W_2^{(n)}(s))|] < \infty$ for $T > 0$ and $h \in \overline{\mathcal{C}}$,
799 and $\mathbb{E} \left[\left| G_n h(W_2^{(n)}(t)) - Gh(W_2(t)) \right| \right] \rightarrow 0$ as $n \rightarrow \infty$ for $h \in \overline{\mathcal{C}}$, where $Gh(W_2(t))$ is the
800 generator of the diffusion process described in Theorem 2.

801 Hereafter, we use the notation $\mathbf{Y}_{1,3:n} = (Y_1, Y_3, \dots, Y_n)$. The latter condition may be reex-
802 pressed as $\mathbb{E}_{\mathbf{X}_{1,2}} [|G_n h(X_2) - Gh(X_2)|] \rightarrow 0$ as $n \rightarrow \infty$, where

$$G_n h(X_2) = n \mathbb{E}_{Y_2} \left[(h(Y_2) - h(X_2)) \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\alpha(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) \right] \right] \quad (\text{A.3})$$

803 with $\alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = 1 \wedge \frac{\pi(\mathbf{Y}^{(n)})}{\pi(\mathbf{x}^{(n)})}$ and π as in (1), and

$$Gh(X_2) = v(\ell, X_1) \left\{ \frac{1}{2} h''(X_2) + \frac{1}{2} \frac{\partial}{\partial X_2} \log f(X_2 | X_1) h'(X_2) \right\}. \quad (\text{A.4})$$

804 There is again a slight abuse of notation as, although h is a function of x_2 only, the generators
805 $G_n h(x_2)$ and $Gh(x_2)$ are functions of x_1, x_2 . Due to the form of (A.4), we can resort to
806 Theorem 2.1 of Chapter 8 in [9] to assert that \mathcal{C}_c^∞ , the space of continuous and infinitely
807 differentiable functions that are compactly supported on \mathbb{R} , forms a core for the generator of
808 the diffusion in Theorem 2. The test function h in (A.3) and (A.4) might then be restricted
809 to functions h belonging to \mathcal{C}_c^∞ .

810 We note that the condition $\sup_n \sup_{s \leq T} \mathbb{E}[|G_n h(W_2^{(n)}(s))|] < \infty$ for $T > 0$ and $h \in \overline{\mathcal{C}}$ may
811 be reexpressed as $\sup_n \mathbb{E}[|G_n h(X_2)|] < \infty$ for $h \in \mathcal{C}_c^\infty$. In fact, it is straight-forward to verify
812 that $\mathbb{E}[(G_n h(X_2))^2] \leq K_h + \mathcal{O}(n^{-1})$ for some $K_h \in (0, \infty)$ which implies that the former is
813 satisfied (this is achieved by considering a function similar to (B.6), in which the acceptance
814 function is Taylor expanded to first order only, and by proceeding as in the proof of Lemma
815 B.3). It also implies the satisfaction of the second applicability condition of Corollary 8.6
816 (Chapter 4 in [9]), which may be simplified as $\limsup_{n \rightarrow \infty} \mathbb{E}[(G_n h(X_2))^2] < \infty$ for $h \in \mathcal{C}_c^\infty$.

817 We now proceed to verify that $G_n h(X_2)$ converges in \mathcal{L}^1 to $Gh(x_2)$. To begin, we have from
818 Lemma B.3 that $\mathbb{E}_{\mathbf{X}_{1,2}} \left[\left| G_n h(X_2) - G_n^{(1)} h(X_2) \right| \right] \rightarrow 0$ as $n \rightarrow \infty$, where $G_n^{(1)} h(x_2)$ is the

819 generator of a diffusive process :

$$\begin{aligned} G_n^{(1)}h(X_2) &= \frac{\ell^2}{2}h''(X_2)\mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\quad + \ell^2 h'(X_2) \frac{\partial}{\partial X_2} \log f(X_2|X_1) \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] . \end{aligned} \quad (\text{A.5})$$

820 Note that it is not necessary to precise this expression further at this stage. Now, we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| G_n^{(1)}h(X_2) - Gh(X_2) \right| \right] \leq \\ &\quad K_1 \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - v(\ell, X_1) \right| \right] \\ &\quad + K_2 \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \frac{\partial}{\partial X_2} \log f(X_2|X_1) \right| \left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - \frac{1}{2}v(\ell, X_1) \right| \right] \end{aligned} \quad (\text{A.6})$$

821 for some $K_1, K_2 > 0$, since $h \in C_c^\infty$ and thus $|h'|$ and $|h''|$ are bounded.

822 Using the triangle inequality, the first term on the RHS of (A.6) satisfies

$$\begin{aligned} &\mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - v(\ell, X_1) \right| \right] \\ &\leq \ell^2 \mathbb{E}_{\mathbf{X}_{1:n}, \mathbf{Y}_{1,3:n}} \left[\left| \alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) - \hat{\alpha}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\quad + \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\hat{\alpha}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - v(\ell, X_1) \right| \right] , \end{aligned}$$

823 where the function $\hat{\alpha}$ is as in Lemma B.4. Using Lemmas B.4 and B.5, the above converges
824 to 0 as $n \rightarrow \infty$. It thus only remains to verify that the second term on the right hand side of
825 (A.6) also converges to 0; Lemma B.6 leads us to that conclusion.

826 B. Appendix : Intermediate results

827 **Lemma B.1.** *As $n \rightarrow \infty$, we have $\mathbb{E} \left[\left| \alpha(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) - \alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right| \right] \rightarrow 0$, with α as in
828 Appendix A.1 and $\alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) = 1 \wedge \exp \left\{ \varepsilon_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right\}$, with*

$$\begin{aligned} \varepsilon_2(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) &= \log \frac{f_1(\mu_n + \frac{\tilde{Y}_1}{\sqrt{n}} | \mathbf{x}_{2:n})}{f_1(\mu_n + \frac{\tilde{x}_1}{\sqrt{n}} | \mathbf{x}_{2:n})} + \sum_{i=2}^n \frac{\partial}{\partial x} \log f(x | \mu_n + \frac{\tilde{Y}_1}{\sqrt{n}}) \Big|_{x=x_i} (Y_i - x_i) \\ &\quad - \frac{\ell^2}{2n} \sum_{i=2}^n \left(\frac{\partial}{\partial x} \log f(x | \mu_n + \frac{\tilde{Y}_1}{\sqrt{n}}) \Big|_{x=x_i} \right)^2 . \end{aligned} \quad (\text{B.1})$$

829 *Proof.* The acceptance function satisfies $\alpha(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) = 1 \wedge \exp \{ \varepsilon(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \}$, where

$$\varepsilon(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) = \log \frac{f_1(\mu_n + \frac{\tilde{Y}_1}{\sqrt{n}}) \prod_{i=2}^n f(x_i | \mu_n + \frac{\tilde{Y}_1}{\sqrt{n}})}{f_1(\mu_n + \frac{\tilde{x}_1}{\sqrt{n}}) \prod_{i=2}^n f(x_i | \mu_n + \frac{\tilde{x}_1}{\sqrt{n}})} + \sum_{i=2}^n \left(\log \frac{f(Y_i | \mu_n + \frac{\tilde{Y}_1}{\sqrt{n}})}{f(x_i | \mu_n + \frac{\tilde{Y}_1}{\sqrt{n}})} \right) .$$

830 Applying obvious changes of variables allows us to express ε in terms of $\mathbf{x}^{(n)}$ and $\mathbf{Y}^{(n)}$:

$$\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = \log \frac{f_1(Y_1|\mathbf{x}_{2:n})}{f_1(x_1|\mathbf{x}_{2:n})} + \sum_{i=2}^n (\log f(Y_i|Y_1) - \log f(x_i|Y_1)).$$

831 Using a second-order Taylor expansion with respect to Y_i around x_i ($i = 2, \dots, n$) to reexpress
832 the last term on the right hand side (RHS) leads to

$$\begin{aligned} \varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) &= \log \frac{f_1(Y_1|\mathbf{x}_{2:n})}{f_1(x_1|\mathbf{x}_{2:n})} + \sum_{i=2}^n \frac{\partial}{\partial x_i} \log f(x_i|Y_1)(Y_i - x_i) \\ &\quad + \frac{1}{2} \sum_{i=2}^n \frac{\partial^2}{\partial U_i^2} \log f(U_i|Y_1)(Y_i - x_i)^2 \end{aligned}$$

833 for some $U_i \in (x_i, Y_i)$ or $U_i \in (Y_i, x_i)$.

834 We note that a candidate Y_1 that does not belong to \mathcal{X}_1 is automatically rejected by the
835 algorithm, *i.e.* functions α , α_2 , α_1 , and α^* are identically 0. Applying changes of variables
836 to the function $\varepsilon_2(\tilde{\mathbf{x}}^{(n)}, \tilde{\mathbf{Y}}^{(n)})$ and using the Lipschitz property of $1 \wedge \exp\{\cdot\}$ along with the
837 fact that $Y_i \sim \mathcal{N}(x_i, \ell^2/n)$, $i = 2, \dots, n$ yield

$$\begin{aligned} \mathbb{E} \left[\left| \alpha(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) - \alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right| \right] &\leq \mathbb{E} \left[\left| \varepsilon(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) - \varepsilon_2(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\leq \mathbb{E} \left[\left| \frac{1}{2} \sum_{i=2}^n \frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1)(Y_i - X_i)^2 + \frac{\ell^2}{2n} \sum_{i=2}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) \right)^2 \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\quad + \frac{\ell^2}{2} \left(\frac{n-1}{n} \right) \mathbb{E} \left[\left| \frac{\partial^2}{\partial U_2^2} \log f(U_2|Y_1) - \frac{\partial^2}{\partial X_2^2} \log f(X_2|Y_1) \right| Z_2^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right], \end{aligned}$$

838 where $Z_2 = \sqrt{n}(Y_2 - X_2)/\ell \sim \mathcal{N}(0, 1)$, and $\mathbb{1}_{\mathcal{X}_1}(y) = 1$ if $y \in \mathcal{X}_1$ and 0 otherwise. From
839 Proposition C.1 in Appendix C, the first term on the RHS converges to 0 as $n \rightarrow \infty$. We now
840 study the second term on the right. Since $Y_2 \rightarrow_{a.s.} x_2$, it implies that $U_2 \rightarrow_{a.s.} x_2$; from the
841 Continuous Mapping Theorem, we have $\left| \frac{\partial^2}{\partial U_2^2} \log f(U_2|Y_1) - \frac{\partial^2}{\partial X_2^2} \log f(X_2|Y_1) \right| \rightarrow_{a.s.} 0$, for all
842 $Y_1 \in \mathcal{X}_1$. Furthermore,

$$\begin{aligned} \mathbb{E} \left[\left(\frac{\partial^2}{\partial U_2^2} \log f(U_2|Y_1) - \frac{\partial^2}{\partial X_2^2} \log f(X_2|Y_1) \right)^2 Z_2^4 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] &\leq 12 \mathbb{E}[K^2(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] \\ &\leq 24 \mathbb{E}[(K(Y_1) - K(X_1))^2 \mathbb{1}_{\mathcal{X}_1}(Y_1)] + 24 \mathbb{E}[K^2(X_1)] \leq 24K^* \frac{\ell^2}{n} + 24 \mathbb{E}[K^2(X_1)] < \infty \end{aligned}$$

843 for some $K^* > 0$ (since $K(x_1)$ satisfies a Lipschitz condition). We conclude, by invoking the
844 Uniform Integrability Theorem, that the second term converges to 0 as $n \rightarrow \infty$.

845 □

846 **Lemma B.2.** *As $n \rightarrow \infty$, we have $\mathbb{E} \left[\left| \alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) - \alpha_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right| \right] \rightarrow 0$, with α_1 as in*
847 *Appendix A.1 and $\alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)})$ as in Lemma B.1.*

848 *Proof.* Applying obvious changes of variables to α_1, α_2 and using the Lipschitz property of
 849 $1 \wedge \exp\{\cdot\}$ yield

$$\begin{aligned} \mathbb{E} \left[\left| \alpha_2(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) - \alpha_1(\tilde{\mathbf{X}}^{(n)}, \tilde{\mathbf{Y}}^{(n)}) \right| \right] &\leq \mathbb{E} \left[\left| \varepsilon_2(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) - \varepsilon_1(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\leq \mathbb{E} \left[\left| \sum_{i=2}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) - \frac{\partial}{\partial X_i} \log f(X_i|X_1) \right) (Y_i - X_i) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \quad (\text{B.2}) \\ &\quad + \frac{\ell^2}{2} \left(\frac{n-1}{n} \right) \mathbb{E} \left[\left| \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) \right)^2 - \left(\frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2 \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right]. \end{aligned}$$

850 The summation in (B.2) is distributed according to a normal random variable with null mean
 851 and variance $\frac{\ell^2}{n} \sum_{i=2}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) - \frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2$. Using Hölder's inequality, the
 852 corresponding expectation is bounded by

$$\left\{ \ell^2 \left(\frac{n-1}{n} \right) \mathbb{E} \left[\left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) - \frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right\}^{1/2}. \quad (\text{B.3})$$

853 Since $Y_1 \rightarrow_{a.s.} x_1$, we use the Continuous Mapping Theorem to affirm that the integrand
 854 converges to 0 almost surely. By assumption, we know that $\mathbb{E}[(\frac{\partial}{\partial X_i} \log f(X_i|X_1))^4] < \infty$.
 855 From the proof of Proposition C.1, we also know that $\mathbb{E}[(\frac{\partial}{\partial X_i} \log f(X_i|Y_1))^4 \mathbb{1}_{\mathcal{X}_1}(Y_1)] < \infty$.
 856 We can thus use the Uniform Integrability Theorem to deduce that the expectation in (B.3)
 857 converges to 0 as $n \rightarrow \infty$. The exact same arguments may be used to conclude that the last
 858 term in (B.2) converges to 0 as $n \rightarrow \infty$. \square

859 **Lemma B.3.** *As $n \rightarrow \infty$ we have $\mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| G_n h(X_2) - G_n^{(1)} h(X_2) \right| \right] \rightarrow 0$, where $G_n h(X_2)$ and
 860 $G_n^{(1)} h(X_2)$ are in (A.3) and (A.5) respectively, with $\mathbf{Y}_{x_2}^{(n)} = (Y_1, x_2, Y_3, \dots, Y_n)$,*

$$g(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = \exp\{\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\} \mathbb{1} \left\{ \exp\{\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\} < 1 \right\}, \quad (\text{B.4})$$

861 *and*

$$\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = \log \frac{f_1(Y_1)}{f_1(x_1)} + \log \frac{f(Y_2|Y_1)}{f(x_2|x_1)} + \sum_{i=3}^n (\log f(Y_i|Y_1) - \log f(x_i|x_1)). \quad (\text{B.5})$$

862 *Proof.* The acceptance rule in (A.3) may be written $\alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = 1 \wedge \exp\{\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\}$,
 863 where the candidates are generated according to $\mathbf{Y}^{(n)} \sim \mathcal{N}(\mathbf{x}^{(n)}, \ell^2 I_n/n)$. We note that a
 864 candidate $Y_1 \notin \mathcal{X}_1$ is automatically rejected by the algorithm, and thus corresponds to an
 865 acceptance probability that is null. It thus not cause any problem to express the acceptance
 866 function as $\alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1)$ wherever necessary.

867 We first Taylor expand the acceptance function $\alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = 1 \wedge \exp\{\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\}$ with
 868 respect to Y_2 around x_2 . As argued in [16], this function is not everywhere differentiable.
 869 However, the points $(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$ at which the derivatives do not exist have a Lebesgue measure
 870 that is either null or converging exponentially to 0 as $n \rightarrow \infty$; hence this shall not cause any
 871 concern when considering expectations of generators. (The latter may happen if f_1 and f
 872 are constant over some interval of the state space, for instance, in which case we could have

873 $\mathbb{P}(\pi(\mathbf{Y}^{(n)}) = \pi(\mathbf{x}^{(n)})) > 0$. The occurrence of such values $\mathbf{x}^{(n)}$ however has a probability
874 converging exponentially rapidly to 0 as $n \rightarrow \infty$.

875 The first-order derivative of $\alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})$ with respect to Y_2 is given by

$$\frac{\partial}{\partial Y_2} \alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = \frac{\partial}{\partial Y_2} \log f(Y_2|Y_1) g(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) ,$$

876 where the function g is as in (B.4); the second-order derivative is expressed as

$$\frac{\partial^2}{\partial Y_2^2} \alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = \left\{ \frac{\partial^2}{\partial Y_2^2} \log f(Y_2|Y_1) + \left(\frac{\partial}{\partial Y_2} \log f(Y_2|Y_1) \right)^2 \right\} g(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) .$$

877 The generator in (A.3) is thus developed as

$$\begin{aligned} G_n h(X_2) &= n \mathbb{E}_{Y_2} [h(Y_2) - h(X_2)] \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &+ n \mathbb{E}_{Y_2} [(h(Y_2) - h(X_2)) (Y_2 - X_2)] \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\frac{\partial}{\partial Y_2} \alpha(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) \Big|_{Y_2=X_2} \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &+ R_n(\mathbf{X}_{1:2}, U_2), \end{aligned} \quad (\text{B.6})$$

878 where

$$R_n(\mathbf{X}_{1:2}, U_2) = \frac{n}{2} \mathbb{E}_{Y_2} \left[(h(Y_2) - h(X_2)) (Y_2 - X_2)^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\frac{\partial^2}{\partial U_2^2} \alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{U_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right] \quad (\text{B.7})$$

879 for some $U_2 \in (X_2, Y_2)$ or $U_2 \in (Y_2, X_2)$. This leads to

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| G_n h(X_2) - G_n^{(1)} h(X_2) \right| \right] &\leq \mathbb{E} [|R_n(\mathbf{X}_{1:2}, U_2)|] \\ &+ \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| n \mathbb{E}_{Y_2} [h(Y_2) - h(X_2)] - \frac{\ell^2}{2} h''(X_2) \right| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right] \quad (\text{B.8}) \\ &+ \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| n \mathbb{E}_{Y_2} [(h(Y_2) - h(X_2)) (Y_2 - X_2)] \right| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\frac{\partial}{\partial Y_2} \alpha(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) \Big|_{Y_2=X_2} \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right. \\ &\quad \left. - \ell^2 h'(X_2) \frac{\partial}{\partial X_2} \log f(X_2|X_1) \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right] . \end{aligned}$$

880 The remainder term in (B.7) converges to 0 in \mathcal{L}^1 , as now detailed. By using a first-order
881 Taylor expansion of h with respect to Y_2 around x_2 along with the fact that $h \in C_c^\infty$, it follows
882 that $|h(Y_2) - h(x_2)| \leq K_1 |Y_2 - x_2|$ for some $K_1 > 0$. Furthermore, since $\frac{\partial}{\partial x_2} \log f(x_2|x_1)$ is
883 Lipschitz continuous on \mathbb{R} for all fixed $x_1 \in \mathcal{X}_1$, then $|\frac{\partial^2}{\partial x_2^2} \log f(x_2|x_1)| \leq K(x_1)$. Using the
884 fact that the function g in (B.4) is bounded by 1, we then write

$$\begin{aligned} \mathbb{E} [|R_n(\mathbf{X}_{1:2}, U_2)|] &\leq \frac{n}{2} K_1 \frac{2^{3/2}}{\sqrt{\pi}} \frac{\ell^3}{n^{3/2}} \mathbb{E}[K(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] \\ &+ \frac{n}{2} K_1 \mathbb{E} \left[|Y_2 - X_2|^3 \left(\frac{\partial}{\partial U_2} \log f(U_2|Y_1) \right)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] . \end{aligned}$$

885 Since $|\frac{\partial}{\partial U_2} \log f(U_2|Y_1)| \leq |\frac{\partial}{\partial x_2} \log f(x_2|x_1)| + L(x_2)|Y_1 - x_1| + K(Y_1)|Y_2 - x_2|$ and $(a+b+c)^2 \leq$
886 $4(a^2 + b^2 + c^2)$ for a, b , and c in \mathbb{R} , then

$$\begin{aligned} \mathbb{E} [|R_n(\mathbf{X}_{1:2}, U_2)|] &\leq \sqrt{\frac{2}{\pi}} K_1 \frac{\ell^3}{n^{1/2}} \left\{ \mathbb{E}[K(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] + 4 \mathbb{E} \left[\left(\frac{\partial}{\partial X_2} \log f(X_2|X_1) \right)^2 \right] \right\} \\ &+ \sqrt{\frac{2}{\pi}} 4 K_1 \frac{\ell^5}{n^{3/2}} \mathbb{E}[L^2(X_2)] + \sqrt{\frac{32}{\pi}} 4 K_1 \frac{\ell^5}{n^{3/2}} \mathbb{E}[K^2(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] . \end{aligned}$$

887 As argued in the proof of Lemma B.1, $\mathbb{E}[K^2(Y_1)\mathbb{1}_{\mathcal{X}_1}(Y_1)] < \infty$; furthermore, the other expectations on the right are finite by assumption. The three terms on the right thus are $\mathcal{O}(n^{-1/2})$,
888 $\mathcal{O}(n^{-3/2})$, and $\mathcal{O}(n^{-3/2})$, which implies that $\mathbb{E}[|R_n(\mathbf{X}_{1:2}, U_2)|] \rightarrow 0$ as $n \rightarrow \infty$.

890 We now turn to the second term on the RHS of (B.8); since the acceptance function takes values in $[0, 1]$, this term is bounded by

$$\mathbb{E}_{X_2} \left[\left| n\mathbb{E}_{Y_2} [h(Y_2) - h(X_2)] - \frac{\ell^2}{2} h''(X_2) \right| \right] \leq \frac{n}{6} \mathbb{E}_{X_2} \left[\left| \mathbb{E}_{Y_2} [h'''(U_2)(Y_2 - X_2)^3] \right| \right]$$

892 for some $U_2 \in (X_2, Y_2)$ or $U_2 \in (Y_2, X_2)$. The term on the right arises from a third-order Taylor expansion of h with respect to Y_2 around X_2 , along with the fact that $Y_2 \sim \mathcal{N}(X_2, \ell^2/n)$.
893 Since $|h'''|$ is bounded by a constant, the previous expression is bounded by $K_2\ell^3/\sqrt{n}$ for
894 some $K_2 > 0$, which converges to 0 as $n \rightarrow \infty$.

896 In a similar fashion, by Taylor expanding h to second order and using the fact that the functions $|h''|$ and g are bounded by $K_3 > 0$ and 1 respectively, the third term on the RHS
897 of (B.8) satisfies

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| n\mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1:n}} \left[(h(Y_2) - h(X_2))(Y_2 - X_2) \frac{\partial}{\partial X_2} \log f(X_2|Y_1) g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right. \right. \\ \left. \left. - \ell^2 h'(X_2) \frac{\partial}{\partial X_2} \log f(X_2|X_1) \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1:3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right| \right] \\ \leq \ell^2 \mathbb{E} \left[|h'(X_2)| \left| \frac{\partial}{\partial X_2} \log f(X_2|Y_1) - \frac{\partial}{\partial X_2} \log f(X_2|X_1) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ + \frac{1}{\sqrt{2\pi}} K_3 \frac{\ell^3}{n^{1/2}} \mathbb{E} \left[\left| \frac{\partial}{\partial X_2} \log f(X_2|Y_1) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right]. \end{aligned}$$

899 From the Lipschitz continuity of $\frac{\partial}{\partial x_2} \log f(x_2|x_1)$ and the fact that h' is bounded in absolute
900 value, the first term on the right of the inequality is bounded by $\ell^2 K_4 \mathbb{E}[L(X_2)|Y_1 - X_1] \leq$
901 $\ell^3 \sqrt{2} K_4 \mathbb{E}[L(X_2)]/\sqrt{\pi n}$ for some $K_4 > 0$; it is thus $\mathcal{O}(n^{-1/2})$. The second term also is
902 $\mathcal{O}(n^{-1/2})$ since $\mathbb{E} \left[\left| \frac{\partial}{\partial X_2} \log f(X_2|Y_1) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] < \infty$ (from the proof of Proposition C.1).

903

□

904 **Lemma B.4.** *As $n \rightarrow \infty$, we have*

$$\mathbb{E}_{\mathbf{X}_{1:n}, \mathbf{Y}_{1:3:n}} \left[\left| \alpha(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) - \hat{\alpha}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \rightarrow 0,$$

905 where $\alpha(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = 1 \wedge \exp\{\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\}$ with ε as in (B.5) and $\hat{\alpha}(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = 1 \wedge$
906 $\exp\{\hat{\varepsilon}(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\}$ with

$$\begin{aligned} \hat{\varepsilon}(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = \log \frac{f_1(Y_1)}{f_1(x_1)} + \log \frac{f(Y_2|Y_1)}{f(x_2|x_1)} + \sum_{i=3}^n \frac{\partial}{\partial x_1} \log f(x_i|x_1)(Y_1 - x_1) \\ + \frac{1}{2} \sum_{i=3}^n \frac{\partial^2}{\partial x_1^2} \log f(x_i|x_1)(Y_1 - x_1)^2 + \sum_{i=3}^n \frac{\partial}{\partial x_i} \log f(x_i|x_1)(Y_i - x_i) - \frac{\ell^2}{2n} \sum_{i=3}^n \left(\frac{\partial}{\partial x_i} \log f(x_i|x_1) \right)^2. \end{aligned} \quad (\text{B.9})$$

907 *Proof.* The function ε in (B.5) is reexpressed as

$$\begin{aligned}\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) &= \log \frac{f_1(Y_1)}{f_1(x_1)} + \log \frac{f(Y_2|Y_1)}{f(x_2|x_1)} + \sum_{i=3}^n (\log f(Y_i|Y_1) - \log f(Y_i|x_1)) \\ &\quad + \sum_{i=3}^n (\log f(Y_i|x_1) - \log f(x_i|x_1)).\end{aligned}$$

908 Using second-order Taylor expansions with respect to Y_i around x_i ($i = 3, \dots, n$) to reexpress
909 the last two terms on the right hand side leads to

$$\begin{aligned}\varepsilon(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) &= \log \frac{f_1(Y_1)}{f_1(x_1)} + \log \frac{f(Y_2|Y_1)}{f(x_2|x_1)} + \sum_{i=3}^n (\log f(x_i|Y_1) - \log f(x_i|x_1)) \\ &\quad + \sum_{i=3}^n \left(\frac{\partial}{\partial x_i} \log f(x_i|Y_1) - \frac{\partial}{\partial x_i} \log f(x_i|x_1) \right) (Y_i - x_i) \\ &\quad + \frac{1}{2} \sum_{i=3}^n \left(\frac{\partial^2}{\partial U_i^2} \log f(U_i|Y_1) - \frac{\partial^2}{\partial U_i^2} \log f(U_i|x_1) \right) (Y_i - x_i)^2 \\ &\quad + \sum_{i=3}^n \frac{\partial}{\partial x_i} \log f(x_i|x_1) (Y_i - x_i) + \frac{1}{2} \sum_{i=3}^n \frac{\partial^2}{\partial x_i^2} \log f(x_i|x_1) (Y_i - x_i)^2 \\ &\quad + \frac{1}{2} \sum_{i=3}^n \left(\frac{\partial^2}{\partial V_i^2} \log f(V_i|x_1) - \frac{\partial^2}{\partial x_i^2} \log f(x_i|x_1) \right) (Y_i - x_i)^2\end{aligned}$$

910 for some $U_i, V_i \in (x_i, Y_i)$ or $U_i, V_i \in (Y_i, x_i)$. Furthermore, by Taylor expanding the third
911 term of the previous expression to second order (with respect to Y_1 around x_1) we obtain

$$\begin{aligned}&\sum_{i=3}^n (\log f(x_i|Y_1) - \log f(x_i|x_1)) \\ &= \sum_{i=3}^n \frac{\partial}{\partial x_1} \log f(x_i|x_1) (Y_1 - x_1) + \frac{1}{2} \sum_{i=3}^n \frac{\partial^2}{\partial x_1^2} \log f(x_i|x_1) (Y_1 - x_1)^2 \\ &\quad + \frac{1}{2} \sum_{i=3}^n \left(\frac{\partial^2}{\partial U_1^2} \log f(x_i|U_1) - \frac{\partial^2}{\partial x_1^2} \log f(x_i|x_1) \right) (Y_1 - x_1)^2\end{aligned}$$

912 for some $U_1 \in (x_1, Y_1)$ or $U_1 \in (Y_1, x_1)$.

913 Using the Lipschitz property of $1 \wedge \exp\{\cdot\}$ yields

$$\begin{aligned}
& \mathbb{E} \left[\left| 1 \wedge \exp\{\varepsilon(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})\} - 1 \wedge \exp\{\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})\} \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\
& \leq \frac{1}{2} \mathbb{E} \left[\left| \sum_{i=3}^n \left(\frac{\partial^2}{\partial U_1^2} \log f(X_i|U_1) - \frac{\partial^2}{\partial X_1^2} \log f(X_i|X_1) \right) (Y_1 - X_1)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right| \right] \\
& \quad + \mathbb{E} \left[\left| \sum_{i=3}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) - \frac{\partial}{\partial X_i} \log f(X_i|X_1) \right) (Y_i - X_i) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right| \right] \\
& \quad + \frac{1}{2} \mathbb{E} \left[\left| \sum_{i=3}^n \left(\frac{\partial^2}{\partial U_i^2} \log f(U_i|Y_1) - \frac{\partial^2}{\partial U_i^2} \log f(U_i|X_1) \right) (Y_i - X_i)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right| \right] \\
& \quad + \mathbb{E} \left[\left| \frac{1}{2} \sum_{i=3}^n \frac{\partial^2}{\partial X_i^2} \log f(X_i|X_1) (Y_i - X_i)^2 + \frac{\ell^2}{2n} \sum_{i=3}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2 \right| \right] \\
& \quad + \frac{1}{2} \mathbb{E} \left[\left| \sum_{i=3}^n \left(\frac{\partial^2}{\partial V_i^2} \log f(V_i|X_1) - \frac{\partial^2}{\partial X_i^2} \log f(X_i|X_1) \right) (Y_i - X_i)^2 \right| \right]. \quad (\text{B.10})
\end{aligned}$$

914 It remains to show that each term on the right hand side converges to 0 as $n \rightarrow \infty$. We look
915 at the first term of (B.10). Using the triangle's inequality and the fact that $(Y_1 - X_1) \sim$
916 $\mathcal{N}(0, \ell^2/n)$, we have

$$\begin{aligned}
& \frac{1}{2} \mathbb{E} \left[\left| \sum_{i=3}^n \left(\frac{\partial^2}{\partial U_1^2} \log f(X_i|U_1) - \frac{\partial^2}{\partial X_1^2} \log f(X_i|X_1) \right) (Y_1 - X_1)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right| \right] \\
& \leq \frac{\ell^2}{2} \left(\frac{n-2}{n} \right) \mathbb{E} \left[\left| \frac{\partial^2}{\partial U_1^2} \log f(X_3|U_1) - \frac{\partial^2}{\partial X_1^2} \log f(X_3|X_1) \right| Z_1^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right],
\end{aligned}$$

917 where $Z_1 = \sqrt{n}(Y_1 - x_1)/\ell \sim \mathcal{N}(0, 1)$. Since $|U_1 - X_1| \leq |Y_1 - X_1|$ and $Y_1 \in \mathcal{X}_1$, then
918 $U_1 \in \mathcal{X}_1$; in addition, $Y_1 \rightarrow_{a.s.} X_1$ implies that $U_1 \rightarrow_{a.s.} X_1$. By the Continuous Mapping
919 Theorem, $\left| \frac{\partial^2}{\partial U_1^2} \log f(X_3|U_1) - \frac{\partial^2}{\partial X_1^2} \log f(X_3|X_1) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \rightarrow_{a.s.} 0$. Since this term is bounded
920 by $2M(X_3) \geq 0$ and that $2\mathbb{E}[M(X_3)Z_1^2] = 2\mathbb{E}[M(X_3)] < \infty$, the Dominated Convergence
921 Theorem can be used to conclude that the first term on the right of (B.10) converges to 0 as
922 $n \rightarrow \infty$.

923 We now consider the second term. Given $x_1, Y_1 \in \mathcal{X}_1$ and $x_i \in \mathbb{R}$ ($i = 3, \dots, n$),

$$\begin{aligned}
& \sum_{i=3}^n \left(\frac{\partial}{\partial x_i} \log f(x_i|Y_1) - \frac{\partial}{\partial x_i} \log f(x_i|x_1) \right) (Y_i - x_i) \\
& \sim \mathcal{N} \left(0, \frac{\ell^2}{n} \sum_{i=3}^n \left(\frac{\partial}{\partial x_i} \log f(x_i|Y_1) - \frac{\partial}{\partial x_i} \log f(x_i|x_1) \right)^2 \right).
\end{aligned}$$

924 Computing the expectation of the half-normal distribution, applying Jensen's inequality (for

925 the square root function, which is concave), and then the triangle inequality lead to

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}_{1,3:n}, \mathbf{Y}_{1,3:n}} \left[\left| \sum_{i=3}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) - \frac{\partial}{\partial X_i} \log f(X_i|X_1) \right) (Y_i - X_i) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right| \right] \\ & \leq \sqrt{\frac{2\ell^2}{\pi} \left(\frac{n-2}{n} \right)} \left(\mathbb{E}_{\mathbf{X}_{1,3}, Y_1} \left[\left(\frac{\partial}{\partial X_3} \log f(X_3|Y_1) - \frac{\partial}{\partial X_3} \log f(X_3|X_1) \right)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right)^{1/2}. \end{aligned}$$

926 Since $Y_1 \xrightarrow{a.s.} X_1$, then $\left(\frac{\partial}{\partial X} \log f(X|Y_1) - \frac{\partial}{\partial X} \log f(X|X_1) \right)^2 \xrightarrow{a.s.} 0$ by the Continuous
927 Mapping Theorem. Furthermore, we know that

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|Y_1) - \frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ & \leq \mathbb{E} [L^4(X)(Y_1 - x_1)^4] = 3 \frac{\ell^4}{n^2} \mathbb{E} [L^4(X)] < \infty ; \end{aligned}$$

928 the Uniform Integrability Theorem can then be used to conclude that the second term on
929 the right hand side of (B.10) converges to 0 as $n \rightarrow \infty$.

930 Using the triangle's inequality and the fact that $(Y_i - X_i) \sim \mathcal{N}(0, \ell^2/n)$ ($i = 3, \dots, n$), the
931 third term on the right hand side of (B.10) is bounded by

$$\frac{\ell^2}{2} \left(\frac{n-2}{n} \right) \mathbb{E} \left[\left| \frac{\partial^2}{\partial U_3^2} \log f(U_3|Y_1) - \frac{\partial^2}{\partial U_3^2} \log f(U_3|X_1) \right| Z_3^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right],$$

932 where $Z_3 = \sqrt{n}(Y_3 - X_3)/\ell \sim \mathcal{N}(0, 1)$. Given that $Y_1 \xrightarrow{a.s.} X_1$, the Continuous Mapping The-
933 orem implies that $\left| \frac{\partial^2}{\partial U_3^2} \log f(U_3|Y_1) - \frac{\partial^2}{\partial U_3^2} \log f(U_3|X_1) \right| \xrightarrow{a.s.} 0$ as $n \rightarrow \infty$. We again invoke
934 the Uniform Integrability Theorem to conclude that the third term on the right converges to
935 0 as $n \rightarrow \infty$, since

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{\partial^2}{\partial U_3^2} \log f(U_3|Y_1) - \frac{\partial^2}{\partial U_3^2} \log f(U_3|X_1) \right)^2 Z_3^4 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ & \leq 6 \mathbb{E} [K^2(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] + 6 \mathbb{E} [K^2(X_1)] < \infty. \end{aligned}$$

936 Replacing Y_1 by X_1 in the proof of Proposition C.1, the fourth term on the right of (B.10) is
937 easily seen to converge towards 0 as $n \rightarrow \infty$. Finally, the last term is bounded by

$$\frac{\ell^2}{2} \left(\frac{n-2}{n} \right) \mathbb{E} \left[\left| \frac{\partial^2}{\partial V_3^2} \log f(V_3|X_1) - \frac{\partial^2}{\partial X_3^2} \log f(X_3|X_1) \right| Z_3^2 \right],$$

938 with $Z_3 = \sqrt{n}(Y_3 - X_3)/\ell$. Given that $Y_3 \xrightarrow{a.s.} X_3$ and $|V_3 - X_3| \leq |Y_3 - X_3|$, we have $V_3 \xrightarrow{a.s.}$
939 X_3 and the Continuous Mapping Theorem implies that the integrand converges to 0 almost
940 surely. Furthermore, the integrand is bounded by $2K(X_1)Z_3^2$ and since $2\mathbb{E}[K(X_1)Z_3^2] =$
941 $2\mathbb{E}[K(X_1)] < \infty$, the Dominated Convergence Theorem is used to conclude the proof.

942 □

943 **Lemma B.5.** *As $n \rightarrow \infty$, we have*

$$\mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\hat{\alpha}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - v(\ell, X_1) \right| \right] \rightarrow 0,$$

944 where $\hat{\alpha}(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)}) = 1 \wedge \exp\{\hat{\varepsilon}(\mathbf{x}^{(n)}, \mathbf{Y}^{(n)})\}$ with the function $\hat{\varepsilon}$ as in (B.9) and the function
945 v as in (6).

946 *Proof.* We have from the triangle inequality

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[\hat{\alpha}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - v(\ell, X_1) \right| \right] \leq \\ \ell^2 \mathbb{E}_{\mathbf{X}_{1:2}, Z_1} \left[\left| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{3:n}} \left[1 \wedge \exp\{\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})\} \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - 2\Phi \left(-\frac{\ell}{2} \gamma^{1/2}(X_1, Z_1) \right) \right| \right], \end{aligned}$$

947 where $Z_1 = \sqrt{n}(Y_1 - x_1)/\ell \sim \mathcal{N}(0, 1)$. From the boundedness of the absolute value in the
948 above expression, it is sufficient to show that, conditionally on $x_1 \in \mathcal{X}_1$, $x_2, Z_1 \in \mathbb{R}$,

$$\left| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{3:n}} \left[1 \wedge \exp\{\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})\} \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - 2\Phi \left(-\frac{\ell}{2} \gamma^{1/2}(X_1, Z_1) \right) \right| \rightarrow 0 \text{ as } n \rightarrow \infty.$$

949 The function $\hat{\varepsilon}$ being evaluated at $Y_2 = x_2$, it is reexpressed as

$$\begin{aligned} \hat{\varepsilon}(\mathbf{x}^{(n)}, (x_1 + \frac{\ell}{\sqrt{n}}Z_1, x_2, \mathbf{Y}_{3:n})) &= \log \frac{f_1(x_1 + \frac{\ell}{\sqrt{n}}Z_1)}{f_1(x_1)} + \log \frac{f(x_2|x_1 + \frac{\ell}{\sqrt{n}}Z_1)}{f(x_2|x_1)} \quad (\text{B.11}) \\ &+ \frac{\ell}{\sqrt{n}} \sum_{i=3}^n \frac{\partial}{\partial x_1} \log f(x_i|x_1) Z_1 + \frac{\ell^2}{2n} \sum_{i=3}^n \frac{\partial^2}{\partial x_1^2} \log f(x_i|x_1) Z_1^2 \\ &+ \sum_{i=3}^n \frac{\partial}{\partial x_i} \log f(x_i|x_1) (Y_i - x_i) - \frac{\ell^2}{2n} \sum_{i=3}^n \left(\frac{\partial}{\partial x_i} \log f(x_i|x_1) \right)^2. \end{aligned}$$

950 In the sequel, we thus condition on $x_1 \in \mathcal{X}_1$, $x_2, Z_1 \in \mathbb{R}$, and study the convergence of
951 every term in (B.11) as $n \rightarrow \infty$. Given any $x_1 \in \mathcal{X}_1$ and $Z_1 \in \mathbb{R}$, $\exists n^* \geq 1$ such that
952 $x_1 + \frac{\ell}{\sqrt{n}}Z_1 \in \mathcal{X}_1$ for all $n \geq n^*$; it therefore follows from the continuity of functions that
953 $\log\{f_1(Y_1)/f_1(x_1)\} \rightarrow 0$ and $\log\{f(x_2|Y_1)/f(x_2|x_1)\} \rightarrow 0$ for any given $x_2 \in \mathbb{R}$. We now show
954 that conditionally on x_1, Z_1 , the remaining terms are asymptotically distributed according
955 to a normal random variable.

956 Given any $x_1 \in \mathcal{X}_1$, $Z_1 \in \mathbb{R}$, applying the Central Limit Theorem to the third term of (B.11)
957 yields

$$\frac{\ell}{\sqrt{n}} Z_1 \sum_{i=3}^n \frac{\partial}{\partial x_1} \log f(X_i|x_1) \rightarrow_d \mathcal{N} \left(0, \ell^2 Z_1^2 \mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right] \right).$$

958 This follows from the regularity assumptions in Section 2, which imply that $\frac{\partial}{\partial x_1} f(x|x_1)$ is
959 locally integrable and thus that we can differentiate outside of the integral sign to obtain

$$\mathbb{E}_X \left[\frac{\partial}{\partial x_1} \log f(X|x_1) \right] = \frac{d}{dx_1} \int_{\mathbb{R}} f(x|x_1) dx = 0.$$

960 To study the fourth term, we condition on $x_1 \in \mathcal{X}_1$, $Z_1 \in \mathbb{R}$ and use the SLLN to get

$$\frac{\ell^2}{2n} Z_1^2 \sum_{i=3}^n \frac{\partial^2}{\partial x_1^2} \log f(X_i|x_1) \rightarrow_{a.s.} \frac{\ell^2}{2} Z_1^2 \mathbb{E}_X \left[\frac{\partial^2}{\partial x_1^2} \log f(X|x_1) \right].$$

961 Again from the regularity assumptions, $\frac{\partial^2}{\partial x_1^2} f(x|x_1)$ is locally integrable and thus the following
962 identity holds :

$$\mathbb{E}_X \left[\frac{\partial^2}{\partial x_1^2} \log f(X|x_1) \right] + \mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right] = \frac{d^2}{dx_1^2} \int_{\mathbb{R}} f(x|x_1) dx = 0;$$

963 therefore, $\mathbb{E}_X \left[\frac{\partial^2}{\partial x_1^2} \log f(X|x_1) \right] = -\mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right]$ for all $x_1 \in \mathcal{X}_1$.

964 Combining the previous developments and making use of Slutsky's Theorem allows us to
965 conclude that given any $x_1 \in \mathcal{X}_1, Z_1 \in \mathbb{R}$,

$$\begin{aligned} & \frac{\ell}{\sqrt{n}} Z_1 \sum_{i=3}^n \frac{\partial}{\partial x_1} \log f(X_i|x_1) + \frac{\ell^2}{2n} Z_1^2 \sum_{i=3}^n \frac{\partial^2}{\partial x_1^2} \log f(X_i|x_1) \\ & \rightarrow_d \mathcal{N} \left(-\frac{\ell^2}{2} Z_1^2 \mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right], \ell^2 Z_1^2 \mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right] \right). \end{aligned}$$

966 Now, given any $x_1 \in \mathcal{X}_1$, the last two terms on the right of (B.11) satisfy

$$\begin{aligned} & \frac{\ell}{\sqrt{n}} \sum_{i=3}^n \frac{\partial}{\partial X_i} \log f(X_i|x_1) Z_i - \frac{\ell^2}{2n} \sum_{i=3}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|x_1) \right)^2, \\ & \rightarrow_p \mathcal{N} \left(-\frac{\ell^2}{2} \mathbb{E}_X \left[\left(\frac{\partial}{\partial X} \log f(X|x_1) \right)^2 \right], \ell^2 \mathbb{E}_X \left[\left(\frac{\partial}{\partial X} \log f(X|x_1) \right)^2 \right] \right); \end{aligned}$$

967 this follows from the WLLN and the fact that $Z_i \sim \mathcal{N}(0, 1)$ independently for $i = 3, \dots, n$.

968 Given $x_1 \in \mathcal{X}_1, Z_1 \in \mathbb{R}$, the two normal random variables just introduced are asymptotically
969 independent (this is easily seen from the fact that $\sqrt{n}(\mathbf{Y}_{3:n} - \mathbf{x}_{3:n})/\ell^2$ is independent of $\mathbf{x}_{3:n}$
970 $\forall n \geq 3$). We therefore conclude that given any $X_1 \in \mathcal{X}_1$ and $Z_1 \in \mathbb{R}$, $\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \rightarrow_d$
971 $\eta(X_1, Z_1)$, where $\eta(x_1, Z_1) \sim \mathcal{N}(-\ell^2 \gamma(x_1, Z_1)/2, \ell^2 \gamma(x_1, Z_1))$, with

$$\gamma(x_1, Z_1) = Z_1^2 \mathbb{E}_X \left[\left(\frac{\partial}{\partial x_1} \log f(X|x_1) \right)^2 \right] + \mathbb{E}_X \left[\left(\frac{\partial}{\partial X} \log f(X|x_1) \right)^2 \right].$$

972 It easily follows from the fact that $\mathbb{1}_{\mathcal{X}_1}(x_1 + \frac{\ell}{\sqrt{n}} Z_1) \rightarrow 1$ given any $x_1 \in \mathcal{X}_1, Z_1 \in \mathbb{R}$, Slutsky's
973 Theorem, and the Continuous Mapping Theorem, that $1 \wedge \exp\{\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})\} \mathbb{1}_{\mathcal{X}_1}(Y_1) \rightarrow_d$
974 $1 \wedge \exp\{\eta(X_1, Z_1)\}$. From Proposition 2.4 in [17], we know that given x_1, Z_1 ,

$$\mathbb{E}_\eta[1 \wedge \exp\{\eta(x_1, Z_1)\}] = 2\Phi \left(-\frac{\ell}{2} \gamma^{1/2}(x_1, Z_1) \right).$$

975 From the convergence in distribution and the boundedness (and thus uniform integrability) of
976 the random variables, the means are known to converge, *i.e.* given any $x_1 \in \mathcal{X}_1$ and $x_2, Z_1 \in \mathbb{R}$

$$\left| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{3:n}} \left[1 \wedge \exp\{\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})\} \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - 2\Phi \left(-\frac{\ell}{2} \gamma^{1/2}(X_1, Z_1) \right) \right| \rightarrow 0,$$

977 which concludes the proof.

978

□

979 **Lemma B.6.** *As $n \rightarrow \infty$, we have*

$$\mathbb{E}_{\mathbf{X}_{1:2}} \left[\left| \frac{\partial}{\partial X_2} \log f(X_2|X_1) \right| \left| \ell^2 \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{1,3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - \frac{1}{2} v(\ell, X_1) \right| \right] \rightarrow 0,$$

980 where the function g is as in (B.4).

981 *Proof.* Making use of the triangle inequality, we may bound the expectation in the statement
 982 of the lemma by

$$\ell^2 \mathbb{E}_{\mathbf{X}_{1:2}, Z_1} \left[\left| \frac{\partial}{\partial X_2} \log f(X_2|X_1) \right| \left| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - \Phi \left(-\frac{\ell}{2} \gamma^{1/2}(X_1, Z_1) \right) \right| \right],$$

983 which is itself bounded by $2\mathbb{E} \left[\left| \frac{\partial}{\partial X_2} \log f(X_2|X_1) \right| \right] < \infty$ since each term in the difference is
 984 bounded by 1 in absolute value. We can thus use the Dominated Convergence Theorem to
 985 bring the limit inside the first expectation. To conclude the proof, all is left to do is to verify
 986 that given any $X_1 \in \mathcal{X}_1$, $X_2, Z_1 \in \mathbb{R}$,

$$\left| \mathbb{E}_{\mathbf{X}_{3:n}, \mathbf{Y}_{3:n}} \left[g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] - \Phi \left(-\frac{\ell}{2} \gamma^{1/2}(X_1, Z_1) \right) \right| \rightarrow 0,$$

987 where $Z_1 = \sqrt{n}(Y_1 - X_1)/\ell$.

988 In the proof of Lemma B.4 we have verified, among other things, that $\mathbb{E}[|\varepsilon(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) -$
 989 $\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})| \mathbb{1}_{\mathcal{X}_1}(Y_1)] \rightarrow 0$ as $n \rightarrow \infty$. This \mathcal{L}^1 -convergence thus entails that $|\varepsilon(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) -$
 990 $\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)})| \mathbb{1}_{\mathcal{X}_1}(Y_1) \rightarrow_p 0$. From the proof of Lemma B.5 we know that given any $X_1 \in$
 991 \mathcal{X}_1 and $X_2, Z_1 \in \mathbb{R}$, $\hat{\varepsilon}(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \rightarrow_d \eta(X_1, Z_1)$. Using Slutsky's Theorem, these
 992 convergences imply that, conditionally on $X_1 \in \mathcal{X}_1$ and $X_2, Z_1 \in \mathbb{R}$, $\varepsilon(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \rightarrow_d$
 993 $\eta(X_1, Z_1)$.

994 From the Continuous Mapping Theorem, we deduce that given any $X_1 \in \mathcal{X}_1$, $Z_1 \in \mathbb{R}$,

$$g(\mathbf{X}^{(n)}, \mathbf{Y}_{X_2}^{(n)}) \mathbb{1}_{\mathcal{X}_1}(Y_1) \rightarrow_d \exp\{\eta(X_1, Z_1)\} \mathbb{1} \{ \exp\{\eta(X_1, Z_1)\} < 1 \}.$$

995 The function under study is obviously not continuous; however, the discontinuities of the
 996 function on the right have null Lebesgue measure and thus the Continuous Mapping Theorem
 997 is applicable as stated in [8] (Theorem 5.1 and its corollaries).

998 By examining the proof of Proposition 2.4 in [17] we obtain, conditionally on $X_1 \in \mathcal{X}_1$,
 999 $Z_1 \in \mathbb{R}$,

$$\mathbb{E}_\eta \left[\exp\{\eta(X_1, Z_1)\} \mathbb{1} \{ \exp\{\eta(X_1, Z_1)\} < 1 \} \right] = \Phi \left(-\frac{\ell}{2} \gamma^{1/2}(X_1, Z_1) \right).$$

1000 From the convergence in distribution and the fact that the random variables under consid-
 1001 eration are bounded (and thus uniformly integrable), the means are known to converge; this
 1002 concludes the proof of the lemma. \square

1003 C. Appendix

1004 **Proposition C.1.** *Define*

$$W(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) = \frac{1}{2} \sum_{i=2}^n \frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1) (Y_i - X_i)^2 + \frac{\ell^2}{2n} \sum_{i=2}^n \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) \right)^2;$$

1005 *then, $\mathbb{E}[|W(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)})| \mathbb{1}_{\mathcal{X}_1}(Y_1)] \rightarrow 0$ as $n \rightarrow \infty$.*

1006 *Proof.* By Jensen's inequality, $\mathbb{E}[|W|] \leq \sqrt{\mathbb{E}[W^2]}$. Developing the square and taking the
 1007 expectation with respect to $\mathbf{Y}_{2:n}$ yield

$$\begin{aligned} \mathbb{E}_{\mathbf{Y}_{2:n}} \left[W^2(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)}) \right] &= \frac{\ell^4}{2n^2} \sum_{i=2}^n \left(\frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1) \right)^2 \\ &\quad + \frac{\ell^4}{4n^2} \left\{ \sum_{i=2}^n \left(\frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1) + \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) \right)^2 \right) \right\}^2, \end{aligned}$$

1008 which implies

$$\begin{aligned} \mathbb{E}_{\mathbf{Y}_{2:n}} \left[|W(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)})| \right] &\leq \frac{\ell^2}{\sqrt{2n}} \left(\frac{1}{n} \sum_{i=2}^n \left(\frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1) \right)^2 \right)^{1/2} \\ &\quad + \frac{\ell^2}{2} \left| \frac{1}{n} \sum_{i=2}^n \left(\frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1) + \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) \right)^2 \right) \right|. \end{aligned}$$

1009 Reapplying Jensen's inequality on the first term and developing the second term lead to

$$\begin{aligned} \mathbb{E}[|W(\mathbf{X}^{(n)}, \mathbf{Y}^{(n)})| \mathbb{1}_{\mathcal{X}_1}(Y_1)] &\leq \frac{\ell^2}{\sqrt{2n}} \left\{ \mathbb{E} \left[\left(\frac{\partial^2}{\partial X^2} \log f(X|Y_1) \right)^2 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \right\}^{1/2} \\ &\quad + \frac{\ell^2}{2} \mathbb{E} \left[\left| \frac{1}{n} \sum_{i=2}^n \left(\frac{\partial^2}{\partial X_i^2} \log f(X_i|X_1) + \left(\frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2 \right) \right| \right] \\ &\quad + \frac{\ell^2}{2} \left(\frac{n-1}{n} \right) \mathbb{E} \left[\left| \frac{\partial^2}{\partial X_i^2} \log f(X_i|Y_1) - \frac{\partial^2}{\partial X_i^2} \log f(X_i|X_1) \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\quad + \frac{\ell^2}{2} \left(\frac{n-1}{n} \right) \mathbb{E} \left[\left| \left(\frac{\partial}{\partial X_i} \log f(X_i|Y_1) \right)^2 - \left(\frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2 \right| \mathbb{1}_{\mathcal{X}_1}(Y_1) \right]. \end{aligned}$$

1010 The first term on the right is bounded by $\ell^2 \{ \mathbb{E} [K^2(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] / (2n) \}^{1/2}$, which converges
 1011 to 0 as $n \rightarrow \infty$ from the argument at the end of the proof of Lemma B.1. From Lemma 12
 1012 in [2], we know that $\frac{\partial}{\partial x} \log f(x|x_1) \rightarrow 0$ as $x \rightarrow \pm\infty$, $\forall x_1 \in \mathcal{X}_1$; hence, given x_1 , we have
 1013 $\mathbb{E}_X \left[\frac{\partial^2}{\partial X^2} \log f(X|x_1) + \left(\frac{\partial}{\partial X} \log f(X|x_1) \right)^2 \right] = \int \frac{\partial^2}{\partial x^2} f(x|x_1) dx = 0$ and by the WLLN,

$$\left| \frac{1}{n} \sum_{i=2}^n \left(\frac{\partial^2}{\partial X_i^2} \log f(X_i|X_1) + \left(\frac{\partial}{\partial X_i} \log f(X_i|X_1) \right)^2 \right) \right| \rightarrow_p 0.$$

1014 To invoke the Uniform Integrability Theorem for the second term, we use the finiteness of
 1015 $\mathbb{E} \left[\left(\frac{\partial^2}{\partial X^2} \log f(X|X_1) \right)^2 \right] \leq \mathbb{E}[K^2(X_1)]$ and $\mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \right]$.

1016 From $Y_1 \rightarrow_{a.s.} x_1$ and the Continuous Mapping Theorem, the integrands of the last two
 1017 terms are seen to converge to 0 almost surely. Since $\mathbb{E}[K^2(Y_1) \mathbb{1}_{\mathcal{X}_1}(Y_1)] < \infty$ (Section A.2)
 1018 and $\mathbb{E}[K^2(X_1)] < \infty$, the third term converges to 0 using the Uniform Integrability Theorem.
 1019 We come to the same conclusion for the fourth term, using $\mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \right] < \infty$ and

$$\begin{aligned} \mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|Y_1) \right)^4 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] &\leq 8 \mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|Y_1) - \frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \mathbb{1}_{\mathcal{X}_1}(Y_1) \right] \\ &\quad + 8 \mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \right] \\ &\leq 24 \frac{\ell^4}{n^2} \mathbb{E} [L^4(X)] + 8 \mathbb{E} \left[\left(\frac{\partial}{\partial X} \log f(X|X_1) \right)^4 \right] < \infty. \end{aligned}$$

1020 □