

Unbiased Markov chain Monte Carlo methods with couplings

Pierre E. Jacob and John O’Leary

Harvard University, Cambridge, USA

and Yves F. Atchadé

Boston University, USA

[Read before The Royal Statistical Society at a meeting organized by the Research Section on Wednesday, December 11th, 2019, Professor A. Doucet in the Chair]

Summary. Markov chain Monte Carlo (MCMC) methods provide consistent approximations of integrals as the number of iterations goes to ∞ . MCMC estimators are generally biased after any fixed number of iterations. We propose to remove this bias by using couplings of Markov chains together with a telescopic sum argument of Glynn and Rhee. The resulting unbiased estimators can be computed independently in parallel. We discuss practical couplings for popular MCMC algorithms. We establish the theoretical validity of the estimators proposed and study their efficiency relative to the underlying MCMC algorithms. Finally, we illustrate the performance and limitations of the method on toy examples, on an Ising model around its critical temperature, on a high dimensional variable-selection problem, and on an approximation of the cut distribution arising in Bayesian inference for models made of multiple modules.

Keywords:

1. Introduction

Markov chain Monte Carlo (MCMC) methods constitute a popular class of algorithms to approximate high dimensional integrals arising in statistics and other fields (Liu, 2008; Robert and Casella, 2004; Brooks *et al.*, 2011; Green *et al.*, 2015). These iterative methods provide estimators that are consistent as the number of iterations grows large but potentially biased for any fixed number of iterations, which discourages the parallel execution of many short chains (Rosenthal, 2000). Consequently, efforts have focused on exploiting parallel processors within each iteration (Tjelmeland, 2004; Brockwell, 2006; Lee *et al.*, 2010; Jacob *et al.*, 2011; Calderhead, 2014; Goudie *et al.*, 2017; Yang *et al.*, 2017) and on the design of parallel chains targeting different distributions (Altekar *et al.*, 2004; Wang *et al.*, 2015; Srivastava *et al.*, 2015). Still, MCMC estimators are ultimately justified by asymptotics in the number of iterations, which is discordant with current trends in computing hardware, characterized by increasing parallelism but stagnating clock speeds.

In this paper we propose a general construction to produce unbiased estimators of integrals with respect to a target probability distribution from MCMC kernels. The lack of bias means that these estimators can be implemented on parallel processors in the framework of Glynn and Heidelberger (1991), without communication between processors. Confidence intervals can be

Address for correspondence: Pierre E. Jacob, Department of Statistics, Harvard University, 1 Oxford Street, Cambridge, MA 02138, USA.
E-mail: pjacob@fas.harvard.edu

constructed with asymptotic guarantees in the number of processors, in contrast with standard MCMC confidence intervals that are justified asymptotically in the number of iterations (e.g. Flegal *et al.* (2008), Gong and Flegal (2016), Atchadé (2016) and Vats *et al.* (2018)). The lack of bias has additional benefits, as discussed in Section 5.5 in which we make use of its interplay with the law of iterated expectations to perform modular inference; see also the discussion in Section 6.

Our contribution follows the path breaking work of Glynn and Rhee (2014), which uses couplings to construct unbiased estimators of integrals with respect to an invariant distribution. They illustrated their construction on Markov chains represented by iterated random functions, leveraging the contraction properties of such functions. Glynn and Rhee (2014) also considered Harris recurrent chains for which an explicit minorization condition holds. Previously, McLeish (2011) employed similar debiasing techniques to obtain ‘nearly unbiased’ estimators from a single MCMC chain. More recently Jacob *et al.* (2019) removed the bias from conditional particle filters (Andrieu *et al.*, 2010) by coupling chains so that they meet in finite time. The present paper brings this type of ‘Rhee–Glynn’ construction to generic MCMC algorithms, with a novel analysis of estimator efficiency and a variety of examples. Our proposed construction involves couplings of MCMC algorithms, which we discuss for generic Metropolis–Hastings and Gibbs samplers.

Couplings have been used to study the convergence properties of MCMC algorithms from both theoretical and practical points of view (e.g. Reutter and Johnson (1995), Johnson (1996), Rosenthal (1997), Johnson (1998, 2013), Neal (1999), Roberts and Rosenthal (2004) and Johndrow and Mattingly (2017)). Couplings also underpin perfect samplers (Propp and Wilson, 1996; Murdoch and Green, 1998; Casella *et al.*, 2001; Flegal and Herbei, 2012; Lee *et al.*, 2014; Huber, 2016). A notable aspect of the approach of Glynn and Rhee (2014) that is preserved in our method is that only two chains must be coupled for the proposed estimator to be unbiased, without further assumptions on the state space or target distribution. Thus the approach applies more broadly than perfect samplers (see Glynn (2016)) while yielding unbiased estimators rather than exact samples. Coupling pairs of Markov chains also forms the basis of the approach of Neal (1999), with a similar motivation for parallel computation. The proposed estimation technique also shares aims with regeneration methods (e.g. Mykland *et al.* (1995) and Brockwell and Kadane (2005)), and we propose a numerical comparison in Section 5.2.

In Section 2 we introduce our estimators and present a coupling of random-walk Metropolis–Hastings (MH) chains as an illustration. In Section 3 we establish the efficiency properties of these estimators, discuss the verification of key assumptions and describe the use of the proposed estimators on parallel processors in light of results from for example Glynn and Heidelberger (1991). In Section 4 we describe how to couple some important MCMC algorithms and illustrate the effect of dimension on algorithms’ performance with a multivariate normal distributions target. Section 5 contains more challenging examples including a multimodal target, a comparison with regeneration methods, sampling problems in large dimensional discrete spaces arising in Bayesian variable selection and Ising models, and an application to modular inference. We discuss our findings in Section 6. Scripts in R (R Core Team, 2015) are available from <https://github.com/pierrejacob/unbiasedmcmc> and supplementary materials are available on line.

2. Unbiased estimation from coupled chains

2.1. Rhee–Glynn estimator

Given a target probability distribution π on a Polish space \mathcal{X} and a measurable real-valued test function h that is integrable with respect to π , we want to estimate the expectation $\mathbb{E}_\pi[h(X)] =$

1 $\int h(x)\pi(dx)$. Let P denote a Markov transition kernel on \mathcal{X} that leaves π invariant, and let
 2 π_0 be some initial probability distribution on \mathcal{X} . Our estimators are based on a coupled pair
 3 of Markov chains $(X_t)_{t \geq 0}$ and $(Y_t)_{t \geq 0}$, which marginally start from π_0 and evolve according
 4 to P . In particular, we suppose that \bar{P} is a transition kernel on the joint space $\mathcal{X} \times \mathcal{X}$ such
 5 that $\bar{P}\{(x, y), A \times \mathcal{X}\} = P(x, A)$ and $\bar{P}\{(x, y), \mathcal{X} \times A\} = P(y, A)$ for any $x, y \in \mathcal{X}$ and any mea-
 6 surable set A . We then construct the coupled Markov chain $(X_t, Y_t)_{t \geq 0}$ as follows. We draw
 7 (X_0, Y_0) such that $X_0 \sim \pi_0$ and $Y_0 \sim \pi_0$. Given (X_0, Y_0) we draw $X_1 \sim P(X_0, \cdot)$, and then for any
 8 $t \geq 1$, given $X_0, (X_1, Y_0), \dots, (X_t, Y_{t-1})$, we draw $(X_{t+1}, Y_t) \sim \bar{P}\{(X_t, Y_{t-1}), \cdot\}$. We consider the
 9 following assumptions.

10 *Assumption 1.* As $t \rightarrow \infty$, $\mathbb{E}[h(X_t)] \rightarrow \mathbb{E}_\pi[h(X)]$. Furthermore, there is an $\eta > 0$ and $D < \infty$
 11 such that $\mathbb{E}[|h(X_t)|^{2+\eta}] \leq D$ for all $t \geq 0$.

12 *Assumption 2.* The chains are such that the meeting time $\tau := \inf\{t \geq 1 : X_t = Y_{t-1}\}$ satisfies
 13 $\mathbb{P}(\tau > t) \leq C\delta^t$ for all $t \geq 0$, for some constants $C < \infty$ and $\delta \in (0, 1)$.

14 *Assumption 3.* The chains stay together after meeting, i.e. $X_t = Y_{t-1}$ for all $t \geq \tau$.

15 By construction, each of the marginal chains $(X_t)_{t \geq 0}$ and $(Y_t)_{t \geq 0}$ has initial distribution π_0
 16 and transition kernel P . Assumption 1 requires these chains to result in a uniformly bounded
 17 $(2 + \eta)$ -moment of h ; more discussion on moments of Markov chains can be found in Tweedie
 18 (1983). Since X_0 and Y_0 may be drawn from any coupling of π_0 with itself, it is possible to set
 19 $X_0 = Y_0$. However, X_1 is then generated from $P(X_0, \cdot)$, so that $X_1 \neq Y_0$ in general. Thus one
 20 cannot force the meeting time to be small by setting $X_0 = Y_0$. Assumption 2 puts a condition
 21 on the coupling that is operated by \bar{P} and would not in general be satisfied for an independent
 22 coupling. Coupled kernels must be carefully designed, using for example common random
 23 numbers and maximal couplings, for assumption 2 to be satisfied. We present a simple case in
 24 Section 2.2 and further examples in Section 4. We stress that the state space is not assumed
 25 to be discrete, and that the constants D and η of assumption 1 and C and δ of assumption 2
 26 do not need to be known to implement the approach proposed. Assumption 3 typically holds
 27 by design; coupled chains that stay identical after meeting were termed ‘faithful’ in Rosenthal
 28 (1997).

29 Under these assumptions we introduce the following motivation for an unbiased estimator of
 30 $\mathbb{E}_\pi[h(X)]$, following Glynn and Rhee (2014). We begin by writing $\mathbb{E}_\pi[h(X)]$ as $\lim_{t \rightarrow \infty} \mathbb{E}[h(X_t)]$.
 31 Then, for any fixed $k \geq 0$,

$$\begin{aligned}
 \mathbb{E}_\pi[h(X)] &= \mathbb{E}[h(X_k)] + \sum_{t=k+1}^{\infty} \{\mathbb{E}[h(X_t)] - \mathbb{E}[h(X_{t-1})]\} \quad (\text{expanding the limit as a telescoping sum}) \\
 &= \mathbb{E}[h(X_k)] + \sum_{t=k+1}^{\infty} \{\mathbb{E}[h(X_t)] - \mathbb{E}[h(Y_{t-1})]\} \quad (\text{since the chains have the same marginals}) \\
 &= \mathbb{E}[h(X_k) + \sum_{t=k+1}^{\infty} \{h(X_t) - h(Y_{t-1})\}] \quad (\text{swapping the expectations and limit}) \\
 &= \mathbb{E}[h(X_k) + \sum_{t=k+1}^{\tau-1} \{h(X_t) - h(Y_{t-1})\}] \quad (\text{by assumption 3})
 \end{aligned}$$

32 We note that the sum in the last equation is 0 if $k + 1 > \tau - 1$. The heuristic argument above sug-
 33 gests that the estimator $H_k(X, Y) = h(X_k) + \sum_{t=k+1}^{\tau-1} \{h(X_t) - h(Y_{t-1})\}$ should have expectation
 34 $\mathbb{E}_\pi[h(X)]$. We observe that this estimator requires $\tau - 1$ calls to \bar{P} and $\max(1, k + 1 - \tau)$ calls to
 35 P ; thus under assumption 2 its cost has a finite expectation.

In Section 3 we establish the validity of the estimator under the three conditions above; this formally justifies the swap of expectation and limit. The estimator can be viewed as a debiased version of $h(X_k)$. Unbiasedness is guaranteed for any choice of $k \geq 0$, but both the cost and the variance of $H_k(X, Y)$ are sensitive to k ; see Section 3.1. Thanks to this unbiasedness property, we can sample $R \in \mathbb{N}$ independent copies of $H_k(X, Y)$ in parallel and average the results to estimate $\mathbb{E}_\pi[h(X)]$ consistently as $R \rightarrow \infty$; we defer further considerations on the use of unbiased estimators on parallel processors to Section 3.3.

Before presenting examples and enhancements to the estimator above, we discuss the relationship between our approach and existing work. There is a rich literature applying forward couplings to study Markov chain convergence (Johnson, 1996, 1998, 2013; Thorisson, 2000; Lindvall, 2002; Rosenthal, 2002; Douc *et al.*, 2004; Nikooienejad *et al.*, 2016), and to obtain new algorithms such as perfect samplers (Huber, 2016) and the methods of Neal (1999) and Neal and Pinto (2001). Our approach is closely related to Glynn and Rhee (2014), who employed pairs of Markov chains to obtain unbiased estimators. The present work combines similar arguments with couplings of MCMC algorithms and proposes further improvements to remove bias at a reduced loss of efficiency.

Indeed Glynn and Rhee (2014) did not apply their methodology to the MCMC setting. They considered chains that are associated with contractive iterated random functions (see also Diaconis and Freedman (1999)), and Harris recurrent chains with an explicit minorization condition. A minorization condition refers to a small set \mathcal{C} , $\lambda > 0$, an integer $m \geq 1$ and a probability measure ν such that, for all $x \in \mathcal{C}$ and some measurable set A , $P^m(x, A) \geq \lambda\nu(A)$. Such a condition is said to be explicit if the set, constant and probability measure are known by the user. Finding explicit small sets that are useful in practice can present a technical challenge, even for MCMC experts (see the discussion and references in Cowles and Rosenthal (1998)). When available, explicit minorization conditions can also be employed to identify regeneration times, yielding estimators that are amenable to parallel computation in the framework of Mykland *et al.* (1995) and Brockwell and Kadane (2005). By contrast Johnson (1996, 1998) and Neal (1999) addressed the question of coupling MCMC algorithms so that pairs of chains meet exactly, without analytical knowledge on the target distribution. The present paper focuses on the use of couplings of this type in the framework of Glynn and Rhee (2014).

2.2. Coupled Metropolis–Hastings example

Before further examination of our estimator and its properties, we present a coupling of MH chains that will typically satisfy assumptions 1–3 in realistic settings; this coupling was proposed in Johnson (1998) as part of a method to diagnose convergence. We postpone discussion of other couplings of MCMC algorithms to Section 4. We recall that each iteration t of the MH algorithm (Hastings, 1970) begins by drawing a proposal X^* from a Markov kernel $q(X_t, \cdot)$, where X_t is the current state. The next state is set to $X_{t+1} = X^*$ if $U \leq \pi(X^*)q(X^*, X_t)/\{\pi(X_t)q(X_t, X^*)\}$, where U denotes a uniform random variable on $[0, 1]$, and $X_{t+1} = X_t$ otherwise.

We define a pair of chains so that each proceeds marginally according to the MH algorithm and jointly so that the chains will meet exactly after a random number of steps. We suppose that the pair of chains are in states X_t and Y_{t-1} , and we consider how to generate X_{t+1} and Y_t so that $\{X_{t+1} = Y_t\}$ might occur.

If $X_t \neq Y_{t-1}$, the event $\{X_{t+1} = Y_t\}$ cannot occur if both chains reject their respective proposals, X^* and Y^* . Meeting will occur if these proposals are identical and if both are accepted. Marginally, the proposals follow $X^*|X_t \sim q(X_t, \cdot)$ and $Y^*|Y_{t-1} \sim q(Y_{t-1}, \cdot)$. If $q(x, x^*)$ can be evaluated for all x and x^* , then we can sample from a maximal coupling between the two

proposal distributions, which is a coupling of $q(X_t, \cdot)$ and $q(Y_{t-1}, \cdot)$ maximizing the probability of the event $\{X^* = Y^*\}$. How to sample from maximal couplings of continuous distributions is described in Thorisson (2000) and in Section 4.1. One can accept or reject the two proposals by using a common uniform random variable U . The chains will stay together after they meet: at each step after meeting, the proposals will be identical with probability 1, and jointly accepted or rejected with a common uniform variable. This coupling requires neither explicit minorization conditions nor contractive properties of a random-function representation of the chain.

2.3. Time-averaged estimator

To motivate our next estimator, we note that we can compute $H_k(X, Y)$ for several values of k from the same realization of the coupled chains, and that the average of these is unbiased as well. For any fixed integer m with $m \geq k$, we can run coupled chains for $\max(m, \tau)$ iterations, compute the estimator $H_l(X, Y)$ for each $l \in \{k, \dots, m\}$ and take the average $H_{k:m}(X, Y) = (m - k + 1)^{-1} \sum_{l=k}^m H_l(X, Y)$, as we summarize in algorithm 1 in Table 1. We refer to $H_{k:m}(X, Y)$ as the *time-averaged estimator*; the estimator $H_k(X, Y)$ is retrieved when $m = k$. Alternatively we could average the estimators $H_l(X, Y)$ by using weights $w_l \in \mathbb{R}$ for $l \in \{k, \dots, m\}$, to obtain $\sum_{l=k}^m w_l H_l(X, Y)$. This will be unbiased if $\sum_{l=k}^m w_l = 1$.

Rearranging terms in $(m - k + 1)^{-1} \sum_{l=k}^m H_l(X, Y)$, we can write the time-averaged estimator as

$$H_{k:m}(X, Y) = \frac{1}{m - k + 1} \sum_{l=k}^m h(X_l) + \sum_{l=k+1}^{\tau-1} \min\left(1, \frac{l-k}{m-k+1}\right) \{h(X_l) - h(Y_{l-1})\}. \quad (2.1)$$

The term $(m - k + 1)^{-1} \sum_{l=k}^m h(X_l)$ corresponds to a standard MCMC average with m total iterations and a burn-in period of $k - 1$ iterations. We can interpret the other term as a bias correction. If $\tau \leq k + 1$, then the correction term equals 0. This provides some intuition for the choice of k and m : large k -values lead to the bias correction being equal to 0 with large probability, and large values of m result in $H_{k:m}(X, Y)$ being similar to an estimator obtained from a long MCMC run. Thus we expect the variance of $H_{k:m}(X, Y)$ to be similar to that of MCMC estimators for appropriate choices of k and m .

The estimator $H_{k:m}(X, Y)$ requires $\tau - 1$ calls to \bar{P} and $\max(1, m + 1 - \tau)$ calls to P , which are overall comparable with m calls to P when m is large. Indeed, for the proposed couplings, calls to \bar{P} are approximately twice as expensive as calls to P . Therefore, the cost of $H_{k:m}(X, Y)$ is comparable with $2(\tau - 1) + \max(1, m + 1 - \tau)$ iterations of the underlying MCMC algorithm. Thus both the variance and the cost of $H_{k:m}(X, Y)$ will approach those of MCMC estimators for large values of k and m . This motivates the use of the estimator $H_{k:m}(X, Y)$ with $m > k$, which

Table 1. Algorithm 1: unbiased ‘time-averaged’ estimator $H_{k:m}(X, Y)$ of $\mathbb{E}_\pi[h(X)]$

<p><i>Step 1:</i> draw X_0 and Y_0 from an initial distribution π_0 and draw $X_1 \sim P(X_0, \cdot)$</p> <p><i>Step 2:</i> set $t = 1$; while $t < \max(m, \tau)$, where $\tau = \inf\{t \geq 1 : X_t = Y_{t-1}\}$,</p> <p style="padding-left: 2em;">(a) draw $(X_{t+1}, Y_t) \sim \bar{P}\{(X_t, Y_{t-1}), \cdot\}$,</p> <p style="padding-left: 2em;">(b) set $t \leftarrow t + 1$</p> <p><i>Step 3:</i> for each $l \in \{k, \dots, m\}$, compute $H_l(X, Y) = h(X_l) + \sum_{t=l+1}^{\tau-1} \{h(X_t) - h(Y_{t-1})\}$;</p> <p>return $H_{k:m}(X, Y) = (m - k + 1)^{-1} \sum_{l=k}^m H_l(X, Y)$, or compute $H_{k:m}(X, Y)$ with equation (2.1)</p>
--

enables us to control the loss of efficiency that is associated with the removal of burn-in bias in contrast with the basic estimator $H_k(X, Y)$ of Section 2.1. We discuss the choice of k and m in further detail in Section 3 and in the subsequent experiments. A variant of estimator (2.1) can be obtained by considering a time lag that is greater than 1 between the two chains $(X_t)_{t \geq 0}$ and $(Y_t)_{t \geq 0}$, with the meeting time defined as the first time t for which $\{X_t = Y_{t-\text{lag}}\}$ occurs. This introduces another tuning parameter but was found to be fruitful in Biswas and Jacob (2019).

We conclude this section with a few remarks on practical implementations. First, the test function h does not have to be specified at run time in algorithm 1 in Table 1. One can store the coupled chains and choose the test function later. Also, one typically resorts to thinning the output of an MCMC sampler if the memory cost of storing chains is prohibitive, or if the cost of evaluating the test function of interest is significant compared with the cost of each MCMC iteration (e.g. Owen (2017)). This is feasible in the framework proposed: one could consider a variation of algorithm 1 where each call to the Markov kernels P and \bar{P} would be replaced by multiple calls to them. We also observe that the estimators proposed can take values outside the range of the test function h ; for instance they can take negative values even if the range of the test function contains only non-negative values.

Finally, we stress the difficulty that is inherent in choosing an initial distribution π_0 . The estimators are unbiased for any choice of π_0 , including point masses, but this choice has an effect on both the computing cost and the variance. There is also a choice about whether to draw X_0 and Y_0 independently from π_0 or not; in our experiments we use independent draws. We shall see in Section 5.1 that unfortunate choices of initial distributions can severely affect the performance of the estimators proposed. This suggests trying more than one choice of initialization, especially in the setting of multimodal targets. Overall the choice of π_0 and its relative importance compared with standard MCMC sampling are open questions.

2.4. Signed measure estimator

We can formulate the proposed estimation procedure in terms of a signed measure $\hat{\pi}$ defined by

$$\hat{\pi} = \frac{1}{m-k+1} \sum_{l=k}^m \delta_{X_l} + \sum_{l=k+1}^{\tau-1} \min\left(1, \frac{l-k}{m-k+1}\right) (\delta_{X_l} - \delta_{Y_{l-1}}), \quad (2.2)$$

which is obtained by replacing test function evaluations by delta masses in equation (2.1), as in Section 4 of Glynn and Rhee (2014). The measure $\hat{\pi}$ is of the form $\hat{\pi} = \sum_{l=1}^N \omega_l \delta_{Z_l}$ where the weights satisfy $\sum_{l=1}^N \omega_l = 1$ and where the atoms (Z_l) are values among the history of the coupled chains. Some of the weights (ω_l) may be negative, making $\hat{\pi}$ a signed empirical measure. In this view the unbiasedness property states that $\mathbb{E}[\sum_{l=1}^N \omega_l h(Z_l)] = \mathbb{E}_\pi[h(X)]$ for a test function h .

We can consider the convergence behaviour of $\hat{\pi}^R = R^{-1} \sum_{r=1}^R \hat{\pi}^{(r)}$ towards π , where $(\hat{\pi}^{(r)})$ for $r \in \{1, \dots, R\}$ are independent replications of $\hat{\pi}$. Glynn and Rhee (2014) obtained a Glivenko–Cantelli result for a similar measure related to their estimator. In the current setting, assume for simplicity that π is univariate or else consider only one of its marginals. To emphasize the importance of the number of replications R , we rewrite the weights and atoms as $\hat{\pi}^R = \sum_{l=1}^{N_R} \omega_l \delta_{Z_l}$. Introduce the function $s \mapsto \hat{F}^R(s) = \sum_{l=1}^{N_R} \omega_l \mathbb{1}(Z_l \leq s)$ on \mathbb{R} . Proposition 2 in Section 3 states that \hat{F}^R converges to F as $R \rightarrow \infty$ uniformly with probability 1, where F is the cumulative distribution function of π .

The function $s \mapsto \hat{F}^R(s)$ is not monotonically increasing because of negative weights among (ω_l) , which motivates the following comments regarding the estimation of quantiles of π . Assume from now on that the pairs (ω_l, Z_l) are ordered such that $Z_l \leq Z_{l+1}$. For any $q \in (0, 1)$ there might be more than one index l such that $\sum_{i=1}^{l-1} \omega_i \leq q$ and $\sum_{i=1}^l \omega_i > q$; the quantile estimate might be

defined as Z_l for any such l . The convergence of \hat{F}^R to F indicates that all such estimates are expected to converge to the q th quantile of π . Therefore the signed measure representation leads to a way of estimating quantiles of the target distribution in a consistent way as $R \rightarrow \infty$. The construction of confidence intervals for these quantiles, perhaps by bootstrapping the R independent copies, stands as an interesting area for future research. Another route to estimate quantiles of π would be to project marginals of $\hat{\pi}^R$ onto the space of probability measures, for instance by using a generalization of the Wasserstein metric to signed measures (Mainini, 2012). One could also estimate F by using isotonic regression (Chatterjee *et al.*, 2015), considering $\hat{F}^R(s)$ for various values s as noisy measurements of $F(s)$.

3. Properties and parallel implementation

The proofs of the results of this section are in the on-line supplementary materials. Our first result establishes the basic validity of the estimators proposed.

Proposition 1. Under assumptions 1–3, for all $k \geq 0$ and $m \geq k$, the estimator $H_{k:m}(X, Y)$ has expectation $\mathbb{E}_\pi[h(X)]$, a finite variance and a finite expected computing time.

A direct consequence of proposition 1 is that an average of R independent copies of $H_{k:m}(X, Y)$ converges to $\mathbb{E}_\pi[h(X)]$ as $R \rightarrow \infty$. We discuss more sophisticated results on unbiased estimators and parallel processing in Section 3.3 and other uses of such estimators in Sections 5.5 and 6. Following Glynn and Rhee (2014), we provide proposition 2 on the signed measure estimator (2.2). We recall that such estimators apply to univariate target distributions or to the marginal distributions of a multivariate target.

Proposition 2. Under assumptions 2 and 3, for all $m \geq k \geq 0$, and assuming that $(X_t)_{t \geq 0}$ converges to π in total variation, introduce the function $s \mapsto \hat{F}^R(s) = \sum_{l=1}^{N_R} \omega_l \mathbb{1}(Z_l \leq s)$, where $(\omega_l, Z_l)_{l=1}^{N_R}$ are weighted atoms obtained from R independent copies of $\hat{\pi}$ in equation (2.2). Denote by F the cumulative distribution function of π . Then

$$\sup_{s \in \mathbb{R}} |\hat{F}^R(s) - F(s)| \xrightarrow{R \rightarrow \infty} 0 \quad \text{almost surely.}$$

Section 3.1 studies the variance and efficiency of $H_{k:m}(X, Y)$, Section 3.2 concerns the verification of assumption 2 by using drift conditions and Section 3.2 discusses estimation on parallel processors in the presence of a budget constraint.

3.1. Variance and efficiency

We consider the effect of k and m on the efficiency of the estimators proposed, which will then suggest guidelines for the choice of these tuning parameters. Estimators $H_{k:m}^{(r)}(X, Y)$, for $r = 1, \dots, R$, can be generated independently and averaged. More estimators can be produced in a given computing budget if each estimator is cheaper to produce. The trade-off can be understood in the framework of Glynn and Whitt (1992); (see also Rhee and Glynn (2012) and Glynn and Rhee (2014)), by defining the asymptotic inefficiency as the product of the variance and expected cost of the estimator. That product is the asymptotic variance of $R^{-1} \sum_{r=1}^R H_{k:m}^{(r)}(X, Y)$ as the computational budget, as opposed to the number of estimators R , goes to ∞ (Glynn and Whitt, 1992). Of primary interest is the comparison of this asymptotic inefficiency with the asymptotic variance of standard MCMC estimators. We start by writing the time-averaged estimator (2.1) as

$$H_{k:m}(X, Y) = \text{MCMC}_{k:m} + \text{BC}_{k:m},$$

where $\text{MCMC}_{k:m}$ is the MCMC average $(m - k + 1)^{-1} \sum_{l=k}^m h(X_l)$ and $\text{BC}_{k:m}$ is the bias correction term. The variance of $H_{k:m}(X, Y)$ can be written

$$\mathbb{V}\{H_{k:m}(X, Y)\} = \mathbb{E}\{[\text{MCMC}_{k:m} - \mathbb{E}_\pi[h(X)]]^2\} + 2\mathbb{E}\{[\text{MCMC}_{k:m} - \mathbb{E}_\pi[h(X)]]\text{BC}_{k:m}\} + \mathbb{E}[\text{BC}_{k:m}^2].$$

Defining the mean-squared error of the MCMC estimator as $\text{MSE}_{k:m} = \mathbb{E}\{[\text{MCMC}_{k:m} - \mathbb{E}_\pi[h(X)]]^2\}$, the Cauchy–Schwarz inequality yields

$$\mathbb{V}\{H_{k:m}(X, Y)\} \leq \text{MSE}_{k:m} + 2\sqrt{\text{MSE}_{k:m}}\sqrt{\mathbb{E}[\text{BC}_{k:m}^2]} + \mathbb{E}[\text{BC}_{k:m}^2]. \quad (3.1)$$

To bound $\mathbb{E}[\text{BC}_{k:m}^2]$, we introduce a geometric drift condition on the Markov kernel P .

Assumption 4. The Markov kernel P is π invariant, φ irreducible and aperiodic, and there is a measurable function $V: \mathcal{X} \rightarrow [1, \infty)$, $\lambda \in (0, 1)$, $b < \infty$, and a small set \mathcal{C} such that, for all $x \in \mathcal{X}$,

$$\int P(x, dy)V(y) \leq \lambda V(x) + b\mathbb{1}(x \in \mathcal{C}).$$

We refer the reader to Meyn and Tweedie (2009) for the definitions and core theoretical tools for working with Markov chains on a general state space: in particular chapter 5 for aperiodicity, φ -irreducibility and small sets, and chapter 15 for geometric drift conditions; see also Roberts and Rosenthal (2004). Geometric drift conditions are known to hold for various MCMC algorithms (e.g. Roberts and Tweedie (1996a, b), Jarner and Hansen (2000), Atchadé (2006), Khare and Hobert (2013), Choi and Hobert (2013) and Pal and Khare (2014)). Assumption 4 often plays a central role in establishing geometric ergodicity (e.g. theorem 9 in Roberts and Rosenthal (2004)). We show next that this assumption enables an informative bound on $\mathbb{E}[\text{BC}_{k:m}^2]$.

Proposition 3. Suppose that assumptions 2–4 hold, with a function V for which the integral $\int V(x)\pi_0(dx)$ is finite. If the function h is such that $\sup_{x \in \mathcal{X}} |h(x)|/V(x)^\beta < \infty$ for some $\beta \in [0, \frac{1}{2})$, then for all $m \geq k \geq 0$ we have

$$\mathbb{E}[\text{BC}_{k:m}^2] \leq \frac{C_{\delta, \beta} \delta_\beta^k}{(m - k + 1)^2},$$

for some constants $C_{\delta, \beta} < \infty$, and $\delta_\beta = \delta^{1-2\beta} \in (0, 1)$, with $\delta \in (0, 1)$ as in assumption 2.

Using proposition 3, inequality (3.1) becomes

$$\mathbb{V}[H_{k:m}(X, Y)] \leq \text{MSE}_{k:m} + 2\sqrt{\text{MSE}_{k:m}} \frac{\sqrt{(C_{\delta, \beta} \delta_\beta^k)}}{m - k + 1} + \frac{C_{\delta, \beta} \delta_\beta^k}{(m - k + 1)^2}. \quad (3.2)$$

The variance of $H_{k:m}(X, Y)$ is thus bounded by the mean-squared error MSE of an MCMC estimator plus additive terms that vanish geometrically in k and polynomially in $m - k$.

To facilitate the comparison between the efficiency of $H_{k:m}(X, Y)$ and that of MCMC estimators, we add simplifying assumptions. First, the rightmost terms of inequality (3.2) decrease geometrically with k , at a rate that is driven by $\delta_\beta = \delta^{1-2\beta}$ where δ is as in assumption 2. This motivates a choice of k depending on the distribution of the meeting time τ . In practice, we can sample independent realizations of the meeting time and choose k such that $\mathbb{P}(\tau > k)$ is small, i.e. we choose k as a large quantile of the meeting times.

Dropping the third term on the right-hand side of inequality (3.2), which is smaller than the second term, assuming that $\text{MSE}_{k:m} > 0$ and that $m > \tau$ with large probability, we obtain the approximate inequality

$$\begin{aligned} \mathbb{E}[2(\tau - 1) + \max(1, m + 1 - \tau)] \mathbb{V}\{H_{k:m}(X, Y)\} &\lesssim \{m + \mathbb{E}(\tau)\} \mathbb{V}\{H_{k:m}(X, Y)\} \\ &\lesssim (m - k + 1) \text{MSE}_{k:m} \left\{ 1 + \frac{k + \mathbb{E}(\tau)}{m - k + 1} \right\} \left[1 + \frac{2}{\sqrt{(m - k + 1)}} \sqrt{\left\{ \frac{C_{\delta, \beta} \delta_{\beta}^k}{(m - k + 1) \text{MSE}_{k:m}} \right\}} \right]. \end{aligned}$$

As k increases we expect $(m - k + 1) \text{MSE}_{k:m}$ to converge to $\mathbb{V}\{(m - k + 1)^{-1/2} \sum_{t=k}^m h(X_t)\}$, where X_k would be distributed according to π . Denote this variance by $V_{k,m}$. The limit of $V_{k,m}$ as $m \rightarrow \infty$ is the asymptotic variance of the MCMC estimator, denoted by V_{∞} . Hence, for k and $m - k$ both large, the loss of efficiency of the method compared with standard MCMC methods is approximately $1 + (k + \mathbb{E}[\tau]) / (m - k)$.

This informal series of approximations suggests that we can retrieve an asymptotic efficiency that is comparable with the underlying MCMC estimators with appropriate choices of k and m that depend on the distribution of the meeting time τ . These choices are thus sensitive to the coupling of the chains and not only to the performance of the underlying MCMC algorithm. Choosing m as a multiple of k , such as $5k$ or $10k$, makes intuitive sense when considering that k/m is the proportion of iterations that are simply discarded in the event that $\tau < k$. In other words, the bias of the MCMC algorithm can be removed at the cost of an increased variance, which can in turn be reduced by choosing sufficiently large values of k and m . This results in a trade-off with the desired level of parallelism: one might prefer to keep k and m small, yielding a suboptimal efficiency for $H_{k:m}(X, Y)$, but enabling more independent copies to be generated in a given computing time.

3.2. Verifying assumption 2

We discuss how assumption 4 on the Markov kernel P can be used to verify assumption 2, on the shape of the meeting time distribution. Informally, assumption 4 guarantees that the bivariate chain $\{(X_t, Y_{t-1}), t \geq 1\}$ visits $\mathcal{C} \times \mathcal{C}$ infinitely often, where \mathcal{C} is a small set. If there is a positive probability of the event $\{X_{t+1} = Y_t\}$ for every t such that $(X_t, Y_{t-1}) \in \mathcal{C} \times \mathcal{C}$, then we expect assumption 2 to hold. The next result formalizes that intuition. The proof is based on a modification of an argument by Douc *et al.* (2004). We introduce $\mathcal{D} = \{(x, y) \in \mathcal{X} \times \mathcal{X} : x = y\}$. Then assumption 3 reads $\bar{P}\{(x, x), \mathcal{D}\} = 1$ for all $x \in \mathcal{X}$.

Proposition 4. Suppose that P satisfies assumption 4 with a small set \mathcal{C} of the form $\mathcal{C} = \{x : V(x) \leq L\}$ where $\lambda + b/(1 + L) < 1$. Suppose also that there exists $\epsilon \in (0, 1)$ such that

$$\inf_{(x, y) \in \mathcal{C} \times \mathcal{C}} \bar{P}\{(x, y), \mathcal{D}\} \geq \epsilon. \quad (3.3)$$

Then there is a finite constant C' and a $\kappa \in (0, 1)$, such that, for all $n \geq 1$,

$$\mathbb{P}(\tau > n) \leq C' \pi_0(V) \kappa^n,$$

where $\pi_0(V) = \int V(x) \pi_0(dx)$. Hence assumption 2 holds as long as $\pi_0(V) < \infty$.

If assumption 4 holds with a small set of the form $\mathcal{C} = \{x : V(x) \leq L\}$ for some $L > 0$, then it holds also for $\mathcal{C} = \{x : V(x) \leq L'\}$ for all $L' \geq L$. In that case we can always choose L sufficiently large that $\lambda + b/(1 + L) < 1$. Hence the main restriction in proposition 4 is the assumption that the small sets in assumption 4 are of the form $\{x : V(x) \leq L\}$, i.e. level sets of V . This is known to be true in some cases. For instance it is known from theorem 2.2 of Roberts and Tweedie (1996b) that, for a large class of MH algorithms, any non-empty compact set is a small set, and therefore for these algorithms it suffices to check that the level sets of the drift function V are compact. Common examples of drift functions include $V(x) = c/\sqrt{\pi(x)}$ (Roberts and Tweedie, 1996b; Jarner and Hansen, 2000; Atchade, 2006), $V(x) = c \exp(b|x|)$ (Roberts and Tweedie, 1996a) or

1 the example in Pal and Khare (2014), which all have compact level sets under mild regularity
2 conditions.

3 The work of Middleton *et al.* (2018) contains results that generalize propositions 3 and 4 to
4 Markov chains satisfying polynomial drift conditions (e.g Andrieu and Vihola, (2015)), leading
5 to polynomial tails for the associated meeting times.

6 3.3. Parallel implementation under budget constraints

7
8 Our main motivation for unbiased estimators comes from parallel processing; see Sections 5.5
9 and 6 for other motivations. Independent unbiased estimators with finite variance can be gener-
10 ated on separate machines and combined into consistent and asymptotically normal estimators.
11 If the number of estimators is prespecified, this follows from the central limit theorem for inde-
12 pendent and identically distributed variables. We might prefer to specify a time budget, and to
13 generate as many estimators as possible within the budget. The lack of bias allows the applica-
14 tion of a variety of results on budget-constrained parallel simulations, which we briefly review
15 here, following Glynn and Heidelberger (1990, 1991).

16 We denote the proposed estimator by H and its expectation, which is the object of interest
17 here, by $\pi(h)$. Generating H takes a random time C . We write $N(t)$ for the number of independent
18 copies of H that can be produced by time t . The sequence $(H_n, C_n)_{n \in \mathbb{N}}$ refers to independent and
19 identically distributed copies of (H, C) , so we can write $N(t) = \sup\{n \geq 0 : C_1 + \dots + C_n \leq t\}$, with
20 $N(t) = 0$ if $t < C_1$. We add the subscript p to refer to objects that are associated with processor
21 $p \in \{1, \dots, P\}$.

22 The first result is that the estimator $\bar{H}_p(t)$, defined for all $1 \leq p \leq P$ as 0 if $N_p(t) = 0$, and
23 by the sample average of $H_{p1}, \dots, H_{pN_p(t)}$ otherwise, is biased: $\mathbb{E}[\bar{H}_p(t)] = \mathbb{E}[H] - \mathbb{E}[H \mathbb{1}(C > t)]$.
24 Corollary 6 of Glynn and Heidelberger (1990) states that, if $\mathbb{E}[H \exp(\alpha C)] < \infty$ for some $\alpha > 0$,
25 then the bias is negligible compared with $\exp(-\alpha t)$ as $t \rightarrow \infty$. By the Cauchy–Schwarz inequality,
26 $\mathbb{E}[H \exp(\alpha C)]^2$ is less than the product of $\mathbb{E}[H^2]$ and $\mathbb{E}[\exp(2\alpha C)]$. In our context, $\mathbb{E}[H^2]$ is finite
27 under proposition 1, and $\mathbb{E}[\exp(2\alpha C)]$ is finite for a range of values of α that depends on the
28 value of δ in assumption 2.

29 We can define an unbiased estimator of $\pi(h)$ with a slight modification of $\bar{H}_p(t)$. For all
30 $p \in \{1, \dots, P\}$, set $\tilde{H}_p(t) = \bar{H}_p(t)$ if $N_p(t) > 0$ and $\tilde{H}_p(t) = H_{p1}$ if $N_p(t) = 0$. With $\tilde{N}_p(t) =$
31 $\max\{1, N_p(t)\}$, then $\tilde{H}_p(t)$ is the sample average of $H_{p1}, \dots, H_{p\tilde{N}_p(t)}$. The computation of $\tilde{N}_p(t)$
32 requires the completion of H_{p1} , and thus we cannot necessarily return $\tilde{H}_p(t)$ at time t , in contrast
33 with $\bar{H}_p(t)$. In contrast, we have $\mathbb{E}[\tilde{H}_p(t)] = \mathbb{E}[H] = \pi(h)$, i.e. the estimator is unbiased, provided
34 that $\mathbb{E}[|H|] < \infty$ (corollary 7 of Glynn and Heidelberger (1990)). We denote the average of $\tilde{H}_p(t)$
35 over P processors by $\tilde{H}(P, t) = P^{-1} \sum_{p=1}^P \tilde{H}_p(t)$, which is unbiased for $\pi(h)$.

36 Asymptotic results on $\tilde{H}(P, t)$ can be found in Glynn and Heidelberger (1991) and are sum-
37 marized below. We first have the consistency results: $\lim_{t \rightarrow \infty} \tilde{H}(P, t) = \lim_{P \rightarrow \infty} \tilde{H}(P, t) = \pi(h)$
38 almost surely for all t and P , and, if $\mathbb{E}[|H|^{1+\delta}] < \infty$ for some $\delta > 0$ and if $\{t_P\}$ is a sequence such
39 that $\lim_{P \rightarrow \infty} t_P = \infty$, then $\tilde{H}(P, t_P)$ converges to $\pi(h)$ in probability as $P \rightarrow \infty$. Next, we can
40 construct confidence intervals for $\pi(h)$ based on $\tilde{H}(P, t)$, following the end of section 3 in Glynn
41 and Heidelberger (1991). Indeed, define

$$42 \hat{\sigma}_1^2(P, t) = \frac{1}{P-1} \sum_{p=1}^P \{\tilde{H}_p(t) - \tilde{H}(P, t)\}^2,$$

$$43 \tilde{\tau}(P, t) = \frac{1}{P} \sum_{p=1}^P \frac{1}{\tilde{N}_p(t)} \sum_{n=1}^{\tilde{N}_p(t)} C_{pn},$$

$$\hat{\sigma}_2^2(P, t) = \tilde{\tau}(P, t) \left\{ \frac{1}{P} \sum_{p=1}^P \frac{1}{\tilde{N}_p(t)} \sum_{n=1}^{\tilde{N}_p(t)} H_{pn}^2 - \tilde{H}(P, t)^2 \right\},$$

where $\tilde{N}_p(t) = \max\{1, N_p(t)\}$. Then we have the three following central limit theorems: for fixed t and $P \rightarrow \infty$,

$$\frac{\sqrt{P}}{\hat{\sigma}_1(P, t)} \{ \tilde{H}(P, t) - \pi(h) \} \rightarrow \mathcal{N}(0, 1); \quad (3.4)$$

for fixed P and $t \rightarrow \infty$,

$$\frac{\sqrt{(Pt)}}{\hat{\sigma}_2(P, t)} \{ \tilde{H}(P, t) - \pi(h) \} \rightarrow \mathcal{N}(0, 1); \quad (3.5)$$

if $t_P \rightarrow \infty$ as $P \rightarrow \infty$,

$$\frac{\sqrt{(Pt_P)}}{\hat{\sigma}_2(P, t_P)} \{ \tilde{H}(P, t_P) - \pi(h) \} \rightarrow \mathcal{N}(0, 1). \quad (3.6)$$

These results require moment conditions such as $\mathbb{E}[\tilde{H}_p(t)^2] < \infty$. The central limit theorem (3.4) will be used to construct confidence intervals in Sections 5.3 and 5.4.

We conclude this section with a remark on the setting where t is fixed and the number of processors P goes to ∞ . There, the time to obtain $\tilde{H}(P, t)$ would typically increase with P . Indeed at least one estimator needs to be completed on each processor for $\tilde{H}(P, t)$ to be available. The completion time behaves as the maximum of independent copies of the cost C . Under assumption 2, the completion time for $\tilde{H}(P, t)$ has expectation behaving as $\log(P)$ when $P \rightarrow \infty$, for fixed t . Other tail assumptions (Middleton *et al.*, 2018) would lead to different behaviour for the completion time that is associated with $\tilde{H}(P, t)$.

4. Couplings of Markov chain Monte Carlo algorithms

We consider couplings of various MCMC algorithms that satisfy assumptions 2 and 3. These couplings are widely applicable and do not require extensive analytical knowledge of the target distribution. We stress that they are not optimal in general, and we expect that other constructions would yield more efficient estimators. We begin in Section 4.1 by reviewing maximal couplings.

4.1. Sampling from maximal couplings

A maximal coupling between two distributions p and q on a space \mathcal{X} is a distribution of a pair of random variables (X, Y) that maximizes $\mathbb{P}(X = Y)$, subject to the marginal constraints $X \sim p$ and $Y \sim q$. We write p and q both for these distributions and for their probability density functions with respect to a common dominating measure, and we refer to the uniform distribution on the interval $[a, b]$ by $\mathcal{U}([a, b])$. A procedure to sample from a maximal coupling is described in algorithm 2 in Table 2 see for example section 4.5 of chapter 1 of Thorisson (2000), and Johnson (1998) where it is termed γ -coupling.

We justify algorithm 2 and compute its cost. Denote by (X, Y) the output of the algorithm. First, X follows p from step 1. To prove that Y follows q , introduce a measurable set A . We write $\mathbb{P}(Y \in A) = \mathbb{P}(Y \in A, \text{step 1}) + \mathbb{P}(Y \in A, \text{step 2})$, where the events $\{\text{step 1}\}$ and $\{\text{step 2}\}$ refer to the algorithm terminating at step 1 or 2. We compute

$$\mathbb{P}(Y \in A, \text{step 1}) = \int_A \int_0^\infty \mathbb{1}\{w \leq q(x)\} \frac{\mathbb{1}\{0 \leq w \leq p(x)\}}{p(x)} p(x) dw dx = \int_A \min\{p(x), q(x)\} dx.$$

Table 2. Algorithm 2: sampling from a maximal coupling of p and q

Step 1: sample $X \sim p$ and $W|X \sim \mathcal{U}\{[0, p(X)]\}$: if $W \leq q(X)$, output (X, X)
 Step 2: otherwise, sample $Y^* \sim q$ and $W^*|Y^* \sim \mathcal{U}\{[0, q(Y^*)]\}$ until
 $W^* > p(Y^*)$, and output (X, Y^*)

We can deduce from this that $\mathbb{P}(\text{step 1}) = \int_{\mathcal{X}} \min\{p(x), q(x)\} dx$. For $\mathbb{P}(Y \in A, \text{step 2})$ to be equal to $\int_A [q(x) - \min\{p(x), q(x)\}] dx$, we need

$$\int_A [q(x) - \min\{p(x), q(x)\}] dx = \mathbb{P}(Y \in A | \text{step 2}) \left[1 - \int_{\mathcal{X}} \min\{p(x), q(x)\} dx \right],$$

and we conclude that the distribution of Y given $\{\text{step 2}\}$ should for all x have a density $\tilde{q}(x)$ equal to $[q(x) - \min\{p(x), q(x)\}] / [1 - \int \min\{p(x'), q(x')\} dx']$. Step 2 is a standard rejection sampler using q as a proposal distribution to target \tilde{q} , which concludes the proof that $Y \sim q$. We also confirm that algorithm 2 maximizes the probability of $\{X = Y\}$. Under the algorithm,

$$\mathbb{P}(X = Y) = \mathbb{P}(\text{step 1}) = \int_{\mathcal{X}} \min\{p(x), q(x)\} dx = 1 - d_{\text{TV}}(p, q),$$

where $d_{\text{TV}}(p, q) = \frac{1}{2} \int_{\mathcal{X}} |p(x) - q(x)| dx$ is the total variation distance. By the coupling inequality (Lindvall, 2002), this proves that the algorithm implements a maximal coupling.

To assess the cost of algorithm 2, note that step 1 costs one draw from p , one evaluation from p and one from q . Each attempt in the rejection sampler of step 2 costs one draw from q , one evaluation from p and one from q . Hereafter we refer to the cost of one draw and two evaluations by ‘one unit’. Observe that the probability of acceptance in step 2 is given by $\mathbb{P}\{W^* \geq p(Y^*)\} = 1 - \int_{\mathcal{X}} \min\{p(y), q(y)\} dy$. Then, the number of attempts in step 2 has a geometric distribution with mean $[1 - \int_{\mathcal{X}} \min\{p(y), q(y)\} dy]^{-1}$, and step 2 itself occurs with probability $1 - \int_{\mathcal{X}} \min\{p(y), q(y)\} dy$. Therefore the overall expected cost is two units. The expectation of the cost is the same for all distributions p and q , whereas the variance of the cost depends on $d_{\text{TV}}(p, q)$, and in fact goes to ∞ as this distance goes to 0.

In algorithm 2, the value of X is not used in the generation of Y^* within step 2. In other words, conditionally on $\{X \neq Y\}$, the two output variables are independent. We might prefer to correlate the outputs in the event $\{X \neq Y\}$, e.g. in random-walk MH steps as in the next section. We describe a maximal coupling presented in Bou-Rabee *et al.* (2018). It applies to distributions p and q on \mathbb{R}^d such that $X \sim p$ can be represented as $X = \mu_1 + \Sigma^{1/2} \dot{X}$, and $Y \sim q$ as $Y = \mu_2 + \Sigma^{1/2} \dot{Y}$, where the pair (\dot{X}, \dot{Y}) follows a coupling of some distribution s with itself. The construction requires that s is spherically symmetrical: $s(x) = s(y)$ for all $x, y \in \mathbb{R}^d$ such that $\|x\| = \|y\|$, where $\|\cdot\|$ denotes the Euclidean norm. For instance, if s is a standard multivariate normal distribution, then $X \sim \mathcal{N}(\mu_1, \Sigma)$ and $Y \sim \mathcal{N}(\mu_2, \Sigma)$.

Let $z = \Sigma^{-1/2}(\mu_1 - \mu_2)$ and $e = z/\|z\|$. We independently draw $\dot{X} \sim s$ and $U \sim \mathcal{U}([0, 1])$ and let

$$\dot{Y} = \begin{cases} \dot{X} + z & \text{if } U \leq \min \left\{ 1, \frac{s(\dot{X} + z)}{s(\dot{X})} \right\}, \\ \dot{X} - 2(e' \dot{X})e & \text{otherwise.} \end{cases}$$

This procedure outputs a pair (\dot{X}, \dot{Y}) that follows a coupling of s with itself. We then define $(X, Y) = (\mu_1 + \Sigma^{1/2} \dot{X}, \mu_2 + \Sigma^{1/2} \dot{Y})$. On the event $\{\dot{Y} = \dot{X} + z\}$, we have $X = Y$. On the event $\{\dot{Y} \neq \dot{X} + z\}$, the vector $\dot{X} - 2(e' \dot{X})e$ is the reflection of \dot{X} through the hyperplane orthogonal to

e that passes through the origin. We show that the output (X, Y) follows a maximal coupling of p and q , which we refer to as a maximal coupling with reflection on the residuals, or a ‘reflection maximal coupling’. First we show that \dot{Y} follows s , closely following the argument in Bou-Rabee *et al.* (2018). For a measurable set B , we compute

$$\mathbb{P}(\dot{Y} \in B) = \int \mathbb{1}_B(x+z) \min\left\{1, \frac{s(x+z)}{s(x)}\right\} s(x) dx + \int \mathbb{1}_B\{x - 2(e'x)e\} \max\left\{0, 1 - \frac{s(x+z)}{s(x)}\right\} s(x) dx.$$

The first integral above becomes $\int \mathbb{1}_B(w) \min\{s(w-z), s(w)\} dw$, after a change of variables $w := x+z$. To simplify the second integral we make the change of variables $w := x - 2(e'x)e$. Since this corresponds to a reflection with respect to a plane orthogonal to e , we have $dw = dx$, and $x = w - 2(e'w)e$; thus

$$\int \mathbb{1}_B\{x - 2(e'x)e\} \max\{0, s(x) - s(x+z)\} dx = \int \mathbb{1}_B(w) \max\{0, s(w) - s(w-z)\} dw,$$

where we have used $s\{w - 2(e'w)e\} = s(w)$ and $s\{w - 2(e'w)e + z\} = s(w-z)$, because $\|w - 2(e'w)e\| = \|w\|$ and $\|w - 2(e'w)e + z\| = \|w - z\|$. Summing the two integrals we obtain $\mathbb{P}(\dot{Y} \in B) = \int_B s(w) dw$, so $\dot{Y} \sim s$.

To verify that the procedure corresponds to a maximal coupling of p and q , we observe that

$$\begin{aligned} \mathbb{P}(X \neq Y) &= \mathbb{P}(\dot{Y} \neq \dot{X} + z) = 1 - \int \min\{s(x), s(x+z)\} dx \\ &= 1 - \int \min\{s\{\Sigma^{-1/2}(\tilde{x} - \mu_1)\}, s\{\Sigma^{-1/2}(\tilde{x} - \mu_2)\}\} |\Sigma^{-1/2}| d\tilde{x}, \end{aligned}$$

with the change of variable $\tilde{x} := \mu_1 + \Sigma^{1/2}x$. This is precisely the total variation distance between p and q , on writing their densities in terms of the density of s . Note that the computational cost that is associated with the above sampling technique is deterministic, in contrast with the cost of algorithm 2.

Finally, for discrete distributions with common finite support, a procedure for sampling from a maximal coupling is described in Section 5.4, with a cost that is also deterministic.

4.2. Metropolis–Hastings steps

In Section 2.2 we described a coupling of MH chains due to Johnson (1998); we summarize the coupled kernel $\tilde{P}\{(X_t, Y_{t-1}), \cdot\}$ in the following procedure.

Step 1: sample $(X^*, Y^*) | (X_t, Y_{t-1})$ from a maximal coupling of $q(X_t, \cdot)$ and $q(Y_{t-1}, \cdot)$.

Step 2: sample $U \sim \mathcal{U}([0, 1])$.

Step 3: if $U \leq \min\{1, \pi(X^*)q(X^*, X_t)/\pi(X_t)q(X_t, X^*)\}$, then $X_{t+1} = X^*$; otherwise $X_{t+1} = X_t$.

Step 4: if $U \leq \min\{1, \pi(Y^*)q(Y^*, Y_{t-1})/\pi(Y_{t-1})q(Y_{t-1}, Y^*)\}$, then $Y_t = Y^*$; otherwise $Y_t = Y_{t-1}$.

Here we address the verification of assumptions 1–3 for this algorithm. Assumption 1 can be verified for MH chains under conditions on the target and the proposal (Nummelin, 2002; Roberts and Rosenthal, 2004). In some settings the explicit drift function that is given in theorem 3.2 of Roberts and Tweedie (1996b) may be used to verify assumption 2 as in Section 3.2. The probability of coupling at the next step given that the chains are in X_t and Y_{t-1} can be controlled as follows. First, the probability of proposing the same value X^* depends on the total variation distance between $q(X_t, \cdot)$ and $q(Y_{t-1}, \cdot)$, which is typically strictly positive if X_t and Y_{t-1} are in bounded subsets of \mathcal{X} . Furthermore, the probability of accepting X^* is often strictly positive

1 on bounded subsets of \mathcal{X} , for instance when $\pi(x) > 0$ for all $x \in \mathcal{X}$. Assumption 3 is satisfied by
 2 design thanks to the use of maximal couplings and common uniform variable U in the above
 3 procedure.

4 Different considerations drive the choice of proposal distribution in standard MCMC sam-
 5 pling and in our proposed estimators. In the case of random-walk proposals with variance Σ ,
 6 larger variances lead to smaller total variation distances between $q(X_t, \cdot)$ and $q(Y_{t-1}, \cdot)$ and thus
 7 larger probabilities of proposing identical values. However, meeting events only occur if propo-
 8 sals are accepted, which is unlikely if Σ is too large. This trade-off could lead to a different
 9 choice of Σ than the optima known for the marginal chains (Roberts *et al.*, 1997), and deserves
 10 further investigation.

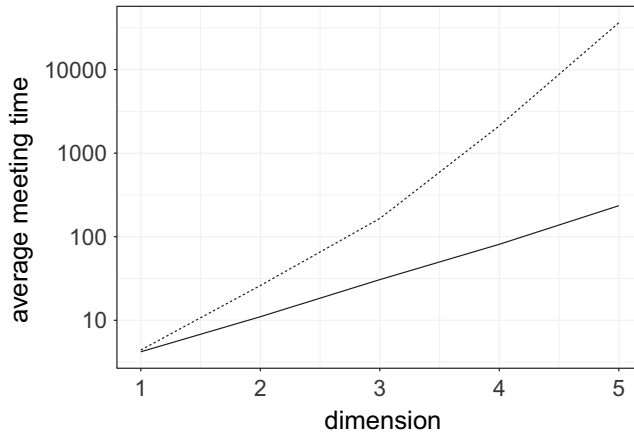
11 We perform experiments with a d -dimensional normal target distribution $\mathcal{N}(0, V)$, where V is
 12 the inverse of a matrix drawn from a Wishart distribution with identity scale matrix and d degrees
 13 of freedom. This setting, which has been borrowed from Hoffman and Gelman (2014), yields
 14 normal distribution targets with strong correlations and a dense precision matrix. Below, each
 15 independent run is performed with an independent draw of V . We consider normal random-walk
 16 proposals with variance Σ set to V/d . The division by d heuristically follows from the scaling
 17 results of Roberts *et al.* (1997). We initialize the chains either from the target distribution, or
 18 from a normal distribution centred at $(1, \dots, 1)$ with identity covariance matrix. We first couple
 19 the proposals with a maximal coupling given by algorithm 2. The resulting average meeting
 20 times, based on 1000 independent runs, are given in Fig. 1(a). The plot indicates an exponential
 21 increase of the average meeting times with the dimension, under both initialization strategies.
 22 In passing, this illustrates that meeting times can be large even if the chains marginally start at
 23 stationarity, i.e. in a setting where there is no burn-in bias.

24 Next we perform the same experiments with the reflection maximum coupling that was de-
 25 scribed in the previous section. The results are shown in Fig. 1(b). The average meeting times
 26 now increase at a rate that appears closer to linear in the dimension. This is to be compared with
 27 established theoretical results on the linear performance of standard MH estimators with respect
 28 to the dimension (Roberts *et al.*, 1997). A formal justification of the scaling that is observed in
 29 Fig. 1(b) is an open question, and so is the design of more effective coupling strategies.

32 4.3. Gibbs sampling

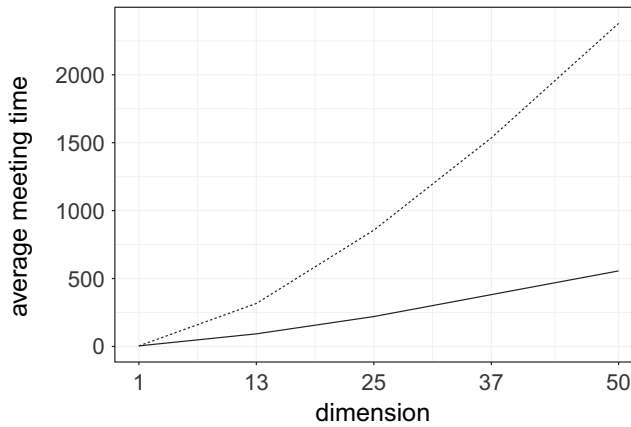
33 Gibbs sampling is another popular class of MCMC algorithms, in which components of a
 34 Markov chain are updated alternately by sampling from the target conditional distributions
 35 (chapter 10 of Robert and Casella (2004)), implemented for example in the software packages
 36 JAGS (Plummer *et al.*, 2003). In Bayesian statistics, these conditional distributions sometimes
 37 belong to a standard family such as the normal, gamma or inverse gamma distributions. Oth-
 38 erwise, the conditional updates might require MH steps. We can introduce couplings in each
 39 conditional update, using either maximal couplings of the target conditionals, if these are stan-
 40 dard distributions, or maximal couplings of the proposal distributions in MH steps targeting
 41 the target conditionals. Controlling the probability of meeting at the next step over a set, as
 42 required for the application of proposition 4, can be done case by case. Drift conditions for
 43 Gibbs samplers also tend to rely on case-by-case arguments (see for example Rosenthal (1996)).

44 Gibbs samplers tend to perform well for targets with weak correlations between the compo-
 45 nents being updated; otherwise Gibbs chains are expected to mix poorly. We perform numerical
 46 experiments on normal target distributions in varying dimensions to observe the effect of cor-
 47 relations on the meeting times of coupled Gibbs chains. For each target $\mathcal{N}(0, V)$, we introduce
 48 an MH-within-Gibbs sampler, where each univariate component i is updated with a single



initialization: — target ··· offset

(a)



initialization: — target ··· offset

(b)

Fig. 1. Scaling of the average meeting time of a coupled MH algorithm with the dimension of the target $\mathcal{N}(0, V)$, where V is the inverse of a Wishart draw, as described in Section 4.2 (the chains are either initialized from the target, or from a normal $\mathcal{N}(\mathbf{1}_d, I_d)$ distribution, where $\mathbf{1}_d$ is a vector of 1s ('offset' in the legend)): (a) using maximal coupling of algorithm 2; (b) using reflection maximal coupling described in Section 4.1

Metropolis step, using normal distribution proposals with variance $V_{i,i}$. Here an iteration of the sampler refers to a complete scan of the components. Fig. 2(a) presents the median meeting times as a function of the dimension, when V is the inverse of a Wishart draw as in the previous section. In this highly correlated setting, the meeting times scale poorly with the dimension. The plot presents the median instead of the average, because we have stopped the runs after 500000 iterations; the median is robust to this truncation, but not the average. We remark that shorter meeting times are obtained when initializing the chains away from the target distribution.

Next we consider a normal distribution target with covariance matrix V defined by $V_{i,j} = 0.5^{-|i-j|}$, which induces weak correlations between components; the inverse of V is tridiagonal. In that case, the same Gibbs sampler performs much more favourably, as we can see from Fig.

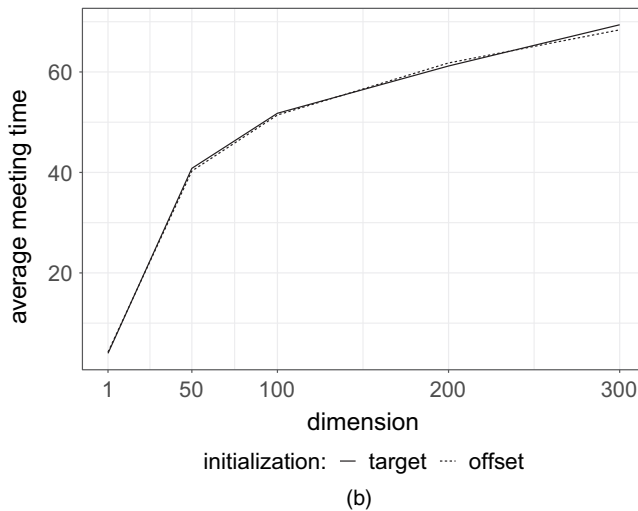
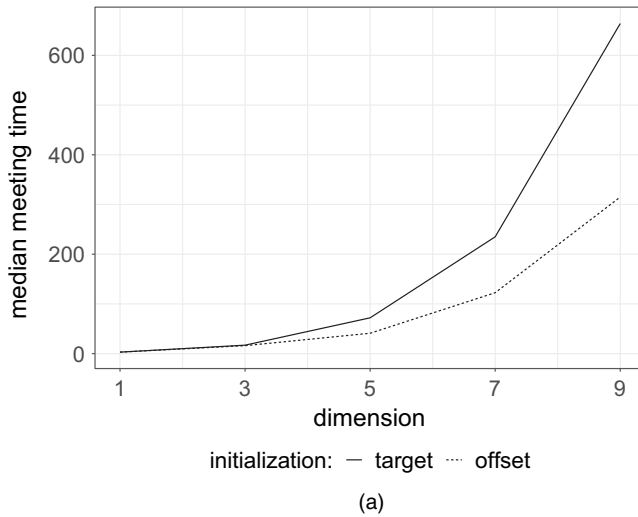


Fig. 2. (a) Scaling of the median and (b) average meeting time of a coupled Gibbs algorithm with the dimension of the target $\mathcal{N}(0, V)$ (the chains are either initialized from the target, or from a normal $\mathcal{N}(\mathbf{1}_d, I_d)$ distribution, where $\mathbf{1}_d$ is a vector of 1s ('offset' in the legend)): (a) the covariance V of the target is generated as the inverse of a Wishart sample, inducing strong correlations; (b) the covariance V is defined as $V_{ij} = 0.5^{-|i-j|}$, inducing weak correlations

2(b). The average meeting times seem to scale sublinearly with the dimension, under both choices of initializations π_0 . Couplings of other Gibbs samplers will be encountered in the numerical experiments of Section 5.

4.4. Coupling of other Markov chain Monte Carlo algorithms

Among extensions of the MH algorithm, Metropolis-adjusted Langevin algorithms (e.g. Roberts and Tweedie, 1996a) are characterized by the use of a proposal distribution given current state X_t that is normal with mean $X_t + h\nabla \log\{\pi(X_t)\}/2$ and variance $h\Sigma$, with tuning parameter $h > 0$ and covariance matrix Σ . Maximal couplings or reflection maximal couplings of the proposals

could be readily implemented to obtain faithful chains. Going further in the use of gradient information, the Hamiltonian or hybrid Monte Carlo algorithm (Duane *et al.*, 1987; Neal, 1993, 2011) is a popular MCMC algorithm for large dimensional targets. In Heng and Jacob (2019), the framework of the present paper is applied to pairs of Hamiltonian Monte Carlo chains, with a focus of the verification of assumptions 1–3 in that context. Such couplings were analysed in detail in Mangoubi and Smith (2017) and Bou-Rabee *et al.* (2018) to obtain convergence rates for the underlying chains. We refer to Heng and Jacob (2019) for more details and provide for completeness some experiments on the normal distribution target that was described above in the on-line supplementary materials.

The present paper generalizes unbiased estimators that are obtained by coupling conditional particle filters in Jacob *et al.* (2019). These algorithms, which were introduced in Andrieu *et al.* (2010), target the distribution of latent processes given observations and fixed parameters for non-linear state space models. The couplings of conditional particle filters in Jacob *et al.* (2019) involve a combination of common random numbers and maximal couplings. Couplings of particle independent MH sampling, which is a particular case of MH sampling with an independent proposal distribution, are simpler to design and are considered in Middleton *et al.* (2019).

The design of generic and efficient MCMC kernels is a topic of active on-going research (see for example Murray *et al.* (2010), Goodman *et al.* (2010), Pollock *et al.* (2016), Vanetti *et al.* (2017) and Titsias and Yau (2017) and references therein). Any new kernel could lead to unbiased estimators with the framework proposed, as long as appropriate couplings can be implemented.

5. Illustrations

Section 5.1 illustrates the effect of k , m and the initial distribution π_0 , identifying a situation where some care is required. Section 5.2 considers the removal of the bias from a Gibbs sampler previously considered for perfect sampling and regeneration methods. Section 5.3 introduces an Ising model and a coupling of a replica exchange algorithm, and we present experiments performed on parallel processors. Section 5.4 considers a high dimensional variable-selection example, with an MH algorithm that has previously been shown to scale linearly with the number of variables. Finally, Section 5.5 focuses on the problem of approximating the cut distribution arising in modular inference, which illustrates the appeal of unbiased estimators beyond parallel computing.

5.1. Bimodal target

We use a bimodal target distribution and a random-walk MH algorithm to illustrate our method and to highlight some of its limitations. In particular, we consider a mixture of univariate normal distributions with density $\pi(x) = 0.5\mathcal{N}(x; -4, 1) + 0.5\mathcal{N}(x; +4, 1)$, which we sample from using random-walk MH with normal proposal distributions of variance $\sigma_q^2 = 9$. This enables regular jumps between the modes of π . We set the initial distribution π_0 to $\mathcal{N}(10, 10^2)$, so that chains are likely to start closer to the mode at 4 than the mode at -4 . Over 1000 independent runs, we find that the meeting time τ has an average of 20 and a 99% quantile of 105.

We consider the task of estimating $\int \mathbb{1}(x > 3)\pi(dx) \approx 0.421$. First, we consider the choice of k and m . Over 1000 independent experiments, we approximate the expected cost $\mathbb{E}[2(\tau - 1) + \max(1, m - \tau + 1)]$ and the variance $\mathbb{V}\{H_{k:m}(X, Y)\}$, and compute the inefficiency as the product of the two (as in Section 3.1). We then divide the inefficiency by the asymptotic variance of the MCMC estimator, which is denoted by V_∞ , which we obtain from 10^6 iterations and a burn-in period of 10^4 by using the R package CODA (Plummer *et al.*, 2006).

1 We present the results in Table 3. First, we see that the inefficiency is sensitive to the choice of
 2 k and m . Second, we see that when k and m are sufficiently large we can retrieve an inefficiency
 3 that is comparable with that of the underlying MCMC algorithm. The ideal choice of k and m
 4 will depend on trade-offs between inefficiency, the desired level of parallelism and the number
 5 of processors that are available. We present a histogram of the target distribution, obtained by
 6 using $k = 200$ and $m = 2000$, in Fig. 3(a). These histograms are produced by averaging unbiased
 7 estimators of expectations of indicator functions, corresponding to consecutive intervals. Con-
 8 fidence intervals at level 95% are obtained from the central limit theorem and are represented
 9 as grey boxes, with vertical bars showing the point estimates.

10 Next, we consider a more challenging case by setting $\sigma_q^2 = 1$, again with $\pi_0 = \mathcal{N}(10, 10^2)$.
 11 These values make it difficult for the chains to jump between the modes of π . Over $R = 1000$
 12 runs we find an average meeting time of 769, with a 99% quantile of 9186. When the chains
 13 start in different modes, the meeting times are often dramatically larger than when the chains
 14 start by the same mode. One can still recover accurate estimates of the target distribution,
 15 but k and m must be set to larger values. With $k = 20000$ and $m = 30000$, we obtain the 95%
 16 confidence interval $[0.397, 0.430]$ for $\int \mathbb{1}(x > 3)\pi(dx) \approx 0.421$. We show a histogram of π in
 17 Fig. 3(b).

Table 3. Cost, variance and inefficiency divided by MCMC asymptotic variance V_∞ , for various choices of k and m , for the test function $h: x \mapsto \mathbb{1}(x > 3)$, in the bimodal target example of Section 5.1

k	m	Cost	Variance	Inefficiency/ V_∞
1	k	37	4.1×10^2	1878.4
1	$10k$	39	3.6×10^2	1703.5
1	$20k$	45	3.0×10^2	1624.8
100	k	119	9.0	130.6
100	$10k$	1019	2.3×10^{-2}	2.9
100	$20k$	2019	7.9×10^{-3}	1.9
200	k	219	2.4×10^{-1}	6.5
200	$10k$	2019	5.3×10^{-3}	1.3
200	$20k$	4019	2.4×10^{-3}	1.2

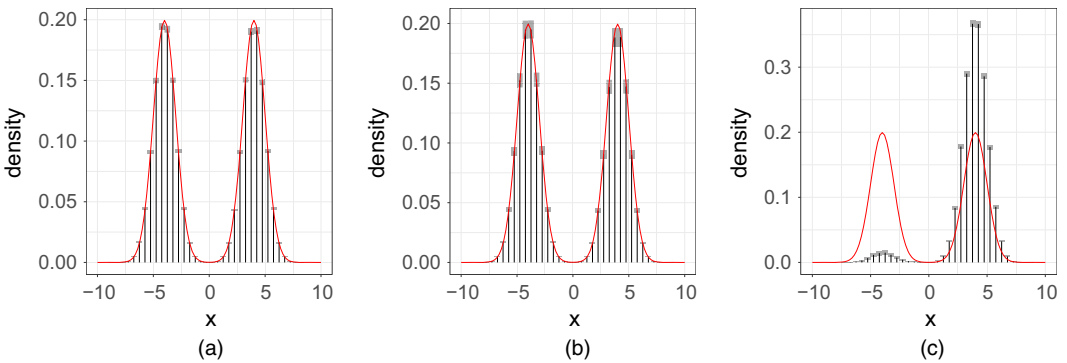


Fig. 3. Histograms of the mixture target distribution of Section 5.1, obtained with the proposed unbiased estimators, based on a normal random-walk MH algorithm, with a proposal variance σ_q^2 and an initial distribution π_0 , over $R = 1000$ experiments (—, target density function): (a) $\sigma_q^2 = 3^2$ and $\pi_0 = \mathcal{N}(10, 10^2)$; (b) $\sigma_q^2 = 1^2$ and $\pi_0 = \mathcal{N}(10, 10^2)$; (c) $\sigma_q^2 = 1^2$ and $\pi_0 = \mathcal{N}(10, 1^2)$

Finally we consider a third case, with $\sigma_q^2 = 1$ as before but now with π_0 set to $\mathcal{N}(10, 1)$. This initialization makes it unlikely for a chain to start near the mode at -4 . The pair of chains typically converge around the mode at 4 and meet in a small number of iterations. Over $R = 1000$ replications, we find an average meeting time of 9 and a 99% quantile of 35. A 95% confidence interval on $\int \mathbb{1}(x > 3)\pi(dx)$ obtained from the estimators with $k = 50$ and $m = 500$ is $[0.799, 0.816]$, which is far from the true value of 0.421. The associated histogram of π is shown in Fig. 3(c).

Sampling 9000 additional estimators yields a 95% confidence interval $[-0.353, 1.595]$, again using $k = 50$ and $m = 500$. Among these extra 9000 values, a few correspond to cases where one chain jumped to the leftmost mode before meeting the other. This resulted in large meeting times and thus a large empirical variance for $H_{k,m}$. On noting a large empirical variance one can then decide to use larger values of k and m . We conclude that, although our estimators are unbiased and are consistent in the limit as $R \rightarrow \infty$, poor performance of the underlying Markov chains combined with ill-chosen initializations can still produce misleading results for any finite R , such as 1000 in this example.

5.2. Gibbs sampler for nuclear pump failure data

Next we consider a classic Gibbs sampler for a model of pump failure counts, which was used for example in Murdoch and Green (1998) to illustrate perfect samplers for continuous distributions, and in Mykland *et al.* (1995) to illustrate their regeneration approach. Here we focus on a comparison with the regeneration approach, which was motivated by similar practical concerns to those in this paper, in particular to avoid an arbitrary choice of burn-in, to construct confidence intervals on the expectations of interest and to make principled use of parallel processors. Mykland *et al.* (1995) showed how to construct regeneration times—random times between which the chain forms independent and identically distributed ‘tours’. They defined a consistent estimator for arbitrary test functions, whose asymptotic variance takes a simple form. The estimator is then obtained by aggregating over these independent tours.

The data consist of operating times $(t_n)_{n=1}^K$ and failure counts $(s_n)_{n=1}^K$ for $K = 10$ pumps at the Farley-1 nuclear power station, as first described in Gaver and O’Muircheartaigh (1987). The model specifies $s_n \sim \text{Poisson}(\lambda_n t_n)$ and $\lambda_n \sim \text{gamma}(\alpha, \beta)$, where $\alpha = 1.802$, $\beta \sim \text{gamma}(\gamma, \delta)$, $\gamma = 0.01$ and $\delta = 1$. The Gibbs sampler for this model consists of the following update steps:

$$\begin{aligned}
 \lambda_n \mid \text{rest} &\sim \text{gamma}(\alpha + s_n, \beta + t_n) && \text{for } n = 1, \dots, K, \\
 \beta \mid \text{rest} &\sim \text{gamma}\left(\gamma + 10\alpha, \delta + \sum_{n=1}^K \lambda_n\right).
 \end{aligned}$$

Here $\text{gamma}(\alpha, \beta)$ refers to the distribution with density $x \mapsto \Gamma(\alpha)^{-1} \beta^\alpha x^{\alpha-1} \exp(-\beta x)$. We initialize all parameter values to 1 (the initialization was not specified in Mykland *et al.* (1995)). To form our estimator we apply maximal couplings at each conditional update of the Gibbs sampler, as described in Section 4.3.

We begin by drawing 1000 meeting times independently. Following the guidelines of Section 3.1, we set $k = 7$, corresponding to the 99% quantile of τ and $m = 10k = 70$. For the regeneration approach, Mykland *et al.* (1995) gave a set of tuning parameters which we adopt below. Applying the regeneration approach to 1000 Gibbs sampler runs of 5000 iterations each, we observe on average 1996 complete tours per run with an average length of 2.50 iterations per tour. These values agree with the count of 1967 tours of average length 2.56 reported in Mykland *et al.* (1995). We observe a posterior mean estimate for β of 2.47 with a variance of 1.89×10^{-4} over the 1000 independent runs, which implies an efficiency value of $(5000 \times 1.89 \times 10^{-4})^{-1} = 1.06$.

1 This exceeds the efficiency of 0.94 that was achieved by our estimator with the choice of $k=7$ and
 2 $m=70$. In contrast, the regeneration approach often requires more extensive analytical work
 3 with the underlying Markov chain; we refer to Mykland *et al.* (1995) for a detailed description.
 4 For reference, the underlying Gibbs sampler achieves an efficiency of 1.08, based on a long
 5 run of 5×10^5 iterations and a burn-in of 10^3 iterations. More extensive comparisons with
 6 other regeneration approaches such as that of Brockwell and Kadane (2005) would deserve
 7 investigation.

10 5.3. Ising model

11 We consider an Ising model on a 32×32 square lattice with periodic boundaries. This provides
 12 a setting where a basic MCMC sampler can mix slowly depending on an inverse temperature
 13 parameter θ , and where a replica exchange strategy as in Geyer (1991) can be helpful. We also
 14 use this example to illustrate the use of our estimators on a large computing cluster, with the
 15 considerations that were reviewed in Section 3.3. For i and j in $\{1, \dots, 32\}^2$ we write $i \sim j$ if
 16 i and j are neighbours in the square lattice with periodic boundaries. We write $x_i \in \{-1, 1\}$
 17 for the spin at location i , and $x = \{x_i\}$ for the full grid. We write $t(x)$ for the 'natural statistic'
 18 $t(x) = 0.5 \sum_{i \in \{1, \dots, 32\}^2} \sum_{j \sim i} x_i x_j$ summing the products of pairs of neighbours. The 0.5-multiplier
 19 here results in each pair of neighbouring sites being counted only once. Under the model, the
 20 probability that is associated with a grid x is $\pi_\theta(x) \propto \exp\{\theta t(x)\}$, where $\theta > 0$ denotes an inverse
 21 temperature parameter that calibrates the degree of correlation between neighbouring sites.

22 We consider a single-site Gibbs sampler, called a heat bath algorithm in this context, to
 23 approximate the distribution π_θ given a value of θ . One iteration of the algorithm consists of
 24 a sweep through all the locations $i \in \{1, \dots, 32\}^2$. For each i we draw x_i from its conditional
 25 distribution under π_θ given all the other spins. It can be checked that the conditional probability
 26 of $\{x_i = 1\}$ given the other spins equals $\exp(\theta s_i) / \{\exp(\theta s_i) + \exp(-\theta s_i)\}$, where s_i denotes the
 27 sum of spins over the four neighbours of i . We initialize the chains by drawing spins uniformly
 28 in $\{-1, 1\}$ at each site, independently across sites.

29 A simple strategy to couple heat bath chains consists of sampling from the maximal coupling
 30 of each conditional distribution. For a grid of θ -values from 0.3 and 0.48, we run 100 pairs of
 31 chains until they meet. We then plot the average meeting time as a function of θ in Fig. 4(a),
 32 noting that the average meeting time increases sharply to values above 10^6 as θ approaches its
 33 critical value (see the related discussion in Propp and Wilson (1996)). We conclude that it would
 34 be expensive to produce unbiased estimators based on the heat bath algorithm for values of θ
 35 above 0.48, for reasons related to the behaviour of the underlying algorithm.

36 There are several ways to address the degeneracy of the heat bath algorithm as θ increases. Spe-
 37 cialized algorithms have been proposed to update groups of spins jointly (Swendsen and Wang,
 38 1987; Wolff, 1989). Here, we consider an approach based on an ensemble of N chains that regu-
 39 larly exchange their states: a technique often termed replica exchange or parallel tempering.
 40 Following for example Geyer (1991), we introduce N chains, $x^{(1)}, \dots, x^{(N)}$, with each $x^{(n)}$ target-
 41 ing $\pi_{\theta^{(n)}}$ with different values of $\theta^{(n)}$ ordered as $\theta^{(1)} < \dots < \theta^{(N)}$. Each iteration of the algorithm
 42 proceeds as follows. With probability $p_{\text{swap}} \in (0, 1)$, for $n \in \{1, \dots, N-1\}$ (sequentially), we pro-
 43 pose exchanging the states $x^{(n)}$ and $x^{(n+1)}$ corresponding to $\theta^{(n)}$ and $\theta^{(n+1)}$. We accept this swap
 44 with probability $\min[1, \pi_{\theta^{(n)}}(x^{(n+1)})\pi_{\theta^{(n+1)}}(x^{(n)}) / \{\pi_{\theta^{(n)}}(x^{(n)})\pi_{\theta^{(n+1)}}(x^{(n+1)})\}]$, which simplifies to
 45 $\min[1, \exp\{(\theta^{(n)} - \theta^{(n+1)})\{t(x^{(n+1)}) - t(x^{(n)})\}\}]$. Otherwise we perform a full sweep of single-site
 46 Gibbs updates, independently across chains.

47 A coupling of this algorithm involves a pair of ensembles with N chains each; the two ensem-
 48 bles are identical if chain n in the first ensemble equals chain n in the second ensemble, for all

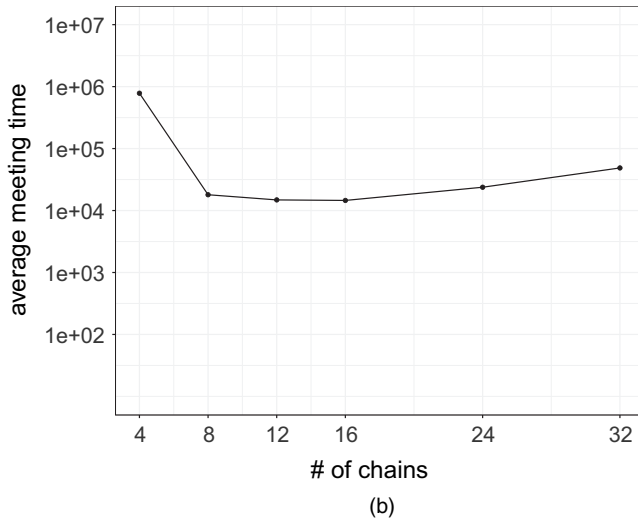
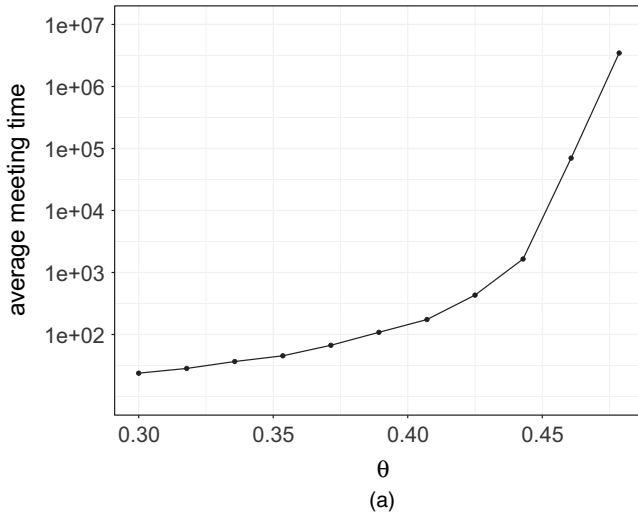


Fig. 4. For the Ising model of Section 5.3 on a 32×32 grid, average meeting times corresponding to different coupled Markov chains: (a) coupled single-site Gibbs sampler, for different inverse temperatures θ ; (b) coupled replica exchange algorithm with N chains, on a grid of values $\theta^{(1)} < \dots < \theta^{(N)}$ with $\theta^{(1)} = 0.3$ and $\theta^{(N)} = 0.55$, for various values of N

$n \in \{1, \dots, N\}$. We use common random numbers to decide whether to perform swap moves or single-site Gibbs moves, and whether to accept the proposed states in the event of a swap move. In the event of a single-site Gibbs move, we maximally couple each conditional update.

Throughout the following experiments we use $p_{\text{swap}} = 0.01$ and introduce an equally spaced grid of θ -values from $\theta^{(1)} = 0.3$ to $\theta^{(N)} = 0.55$ for several choices of N . We note that these grids includes θ -values at which we have seen that the single-site Gibbs sampler mixes poorly. Fig. 4(b) shows the resulting average meeting times over 100 independent runs, as a function of the number of chains N . The average meeting time first decreases with the number of chains but then increases again. A possible explanation is that the mixing of the chains first improves as N increases and then stabilizes; however, it becomes more difficult for the ensembles to meet when

1 N increases since all chains in the ensembles must meet. The minimum average meeting time is
 2 here attained for $N = 16$ chains per ensemble.

3 Setting $N = 16$, $k = 10^5$ and $m = 2 \times 10^5$ we now illustrate the use of the proposed unbiased
 4 estimators on a cluster. The test function is taken as $x \mapsto t(x)$ defined above, so we estimate
 5 $\Sigma_x t(x) \pi_\theta(x)$ for various values of θ . We use 500 processors to generate unbiased estimates with
 6 a time budget of 30 min. Within that time, each processor generated between one and seven
 7 estimators, with an average of 3.7 estimators per processor and a total of 1858 estimators. The
 8 chronology of the generation of these estimates is illustrated in Fig. 5(a). For each processor,
 9 horizontal segments of different colours indicate the duration that is associated with each esti-
 10 mator. The final estimates with standard errors are shown in Fig. 5(b), where we can see that
 11 the standard errors are very small compared with the values of the estimates, for each value of
 12 θ . These standard errors were computed as $\hat{\sigma}_1(P, t)/\sqrt{P}$ following equation (3.4), the central
 13 limit theorem corresponding to the large processor count limit.

14 5.4. Variable selection

15 For our next example we consider a variable-selection problem following Yang *et al.* (2016) to
 16 illustrate the scaling of our proposed method on high dimensional discrete state spaces. For in-
 17 tegers p and n , let $Y \in \mathbb{R}^n$ represent a response variable depending on covariates $X_1, \dots, X_p \in \mathbb{R}^n$.
 18 We consider the task of inferring a binary vector $\gamma \in \{0, 1\}^p$ representing which covariates to
 19 select as predictors of Y , with the convention that X_i is selected if $\gamma_i = 1$. For any γ , we write
 20 $|\gamma| = \sum_{i=1}^p \gamma_i$ for the number of selected covariates and X_γ for the $n \times |\gamma|$ matrix of covariates
 21 chosen by γ . Inference on γ proceeds by way of a linear regression model relating Y to X_γ ,
 22 namely $Y = X_\gamma \beta_\gamma + w$ with $w \sim \mathcal{N}(0, \sigma^2 I_n)$.

23 We assume a prior on γ of $\pi(\gamma) \propto p^{-\kappa|\gamma|} \mathbb{1}(|\gamma| \leq s_0)$. This distribution puts mass only on vectors
 24 γ with fewer than s_0 1s, imposing a degree of sparsity. Given γ we assume a normal prior for the
 25 regression coefficient vector $\beta_\gamma \in \mathbb{R}^{|\gamma|}$ with zero mean and variance $g\sigma^2(X'_\gamma X_\gamma)^{-1}$. Finally, we
 26 give the precision σ^{-2} an improper prior $\pi(\sigma^{-2}) \propto 1/\sigma^{-2}$. This leads to the marginal likelihood

$$27 \pi(Y|X, \gamma) \propto \frac{(1+g)^{-|\gamma|/2}}{\{1+g(1-R_\gamma^2)\}^{n/2}}, \quad R_\gamma^2 = \frac{Y' X_\gamma (X'_\gamma X_\gamma)^{-1} X'_\gamma Y}{Y' Y}.$$

28 To approximate the distribution $\pi(\gamma|X, Y)$, Yang *et al.* (2016) employed an MCMC algorithm
 29 whose kernel P is a mixture of two Metropolis kernels. The first component $P_1(\gamma, \cdot)$ selects a
 30 co-ordinate $i \in \{1, \dots, p\}$ uniformly at random and flips γ_i to $1 - \gamma_i$. The resulting vector γ^*
 31 is then accepted with probability $1 \wedge \pi(\gamma^*|X, Y)/\pi(\gamma|X, Y)$, where $a \wedge b$ denotes $\min(a, b)$ for
 32 $a, b \in \mathbb{R}$. Sampling a vector γ' from the second kernel $P_2(\gamma, \cdot)$ proceeds as follows. If $|\gamma|$ equals 0
 33 or p , then γ' is set to γ . Otherwise, co-ordinates i_0 and i_1 are drawn uniformly among $\{j: \gamma_j = 0\}$
 34 and $\{j: \gamma_j = 1\}$ respectively. The proposal γ^* has $\gamma_{i_0}^* = \gamma_{i_1}$, $\gamma_{i_1}^* = \gamma_{i_0}$ and $\gamma_j^* = \gamma_j$ for the other
 35 components. Then γ' is set to γ^* with probability $1 \wedge \pi(\gamma^*|X, Y)/\pi(\gamma|X, Y)$, and to γ otherwise.
 36 The MCMC kernel $P(\gamma, \cdot)$ targets $\pi(\gamma|X, Y)$ by sampling from $P_1(\gamma, \cdot)$ or from $P_2(\gamma, \cdot)$ with
 37 equal probability. Note that each MCMC iteration can only benefit from parallel processors
 38 to a limited extent, since $|\gamma|$ is always less than s_0 , itself chosen to be a small value; thus the
 39 calculation of R_γ^2 involves only linear algebra of small matrices.

40 We consider the following strategy to couple the above MCMC algorithm. To sample a pair
 41 of states $(\gamma', \tilde{\gamma}')$ given $(\gamma, \tilde{\gamma})$, we first use a common uniform random variable to decide whether
 42 to sample from a coupling \bar{P}_1 of P_1 to itself or a coupling \bar{P}_2 of P_2 to itself. The coupled kernel
 43 $\bar{P}_1\{(\gamma, \tilde{\gamma}), \cdot\}$ proposes flipping the same co-ordinate for both vectors γ and $\tilde{\gamma}$ and then uses a
 44 common uniform random variable in the acceptance step. For the coupled kernel $\bar{P}_2\{(\gamma, \tilde{\gamma}), \cdot\}$,

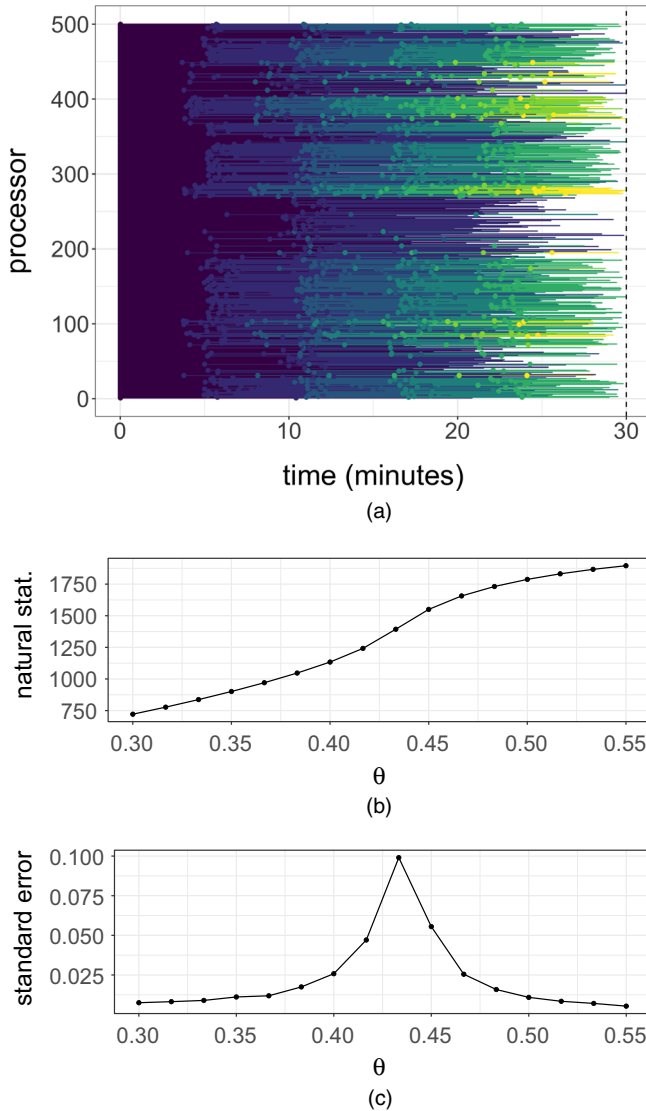


Fig. 5. For the Ising model of Section 5.3 on a 32×32 grid, (a) chronology of the generation of unbiased estimators on 500 processors over 30 min, and (b) estimates of the expected natural statistic $\Sigma_X t(x) \pi_\theta(x)$, and (c) standard errors, for a grid of 16 values $0.3 = \theta^{(1)} < \dots < \theta^{(N)} = 0.55$: results obtained by coupling a replica exchange algorithm with 16 chains

we need to select two pairs of indices: (i_0, \tilde{i}_0) and (i_1, \tilde{i}_1) . We obtain the first pair by sampling from a maximal coupling of the discrete uniform distributions on $\{j: \gamma_j = 0\}$ and $\{j: \tilde{\gamma}_j = 0\}$. This yields indices (i_0, \tilde{i}_0) with the greatest possible probability that $i_0 = \tilde{i}_0$. We use the same approach to sample a pair (i_1, \tilde{i}_1) to maximize the probability that $i_1 = \tilde{i}_1$. Finally we use a common uniform variable to accept or reject the proposals. If either vector γ or $\tilde{\gamma}$ has no 0s or no 1s, then it is kept unchanged.

We recall that one can sample from a maximal coupling of two discrete probability distributions $q = (q_1, \dots, q_N)$ and $\tilde{q} = (\tilde{q}_1, \dots, \tilde{q}_N)$ as follows. First, let $c = (c_1, \dots, c_N)$ be the distribution with probabilities $c_n = (q_n \wedge \tilde{q}_n) / \alpha$ for $\alpha = \sum_{n=1}^N q_n \wedge \tilde{q}_n$ and define residual distributions q' and

\tilde{q}' with probabilities $q'_n = (q_n - \alpha c_n)/(1 - \alpha)$ and $\tilde{q}'_n = (\tilde{q}_n - \alpha c_n)/(1 - \alpha)$. Then, with probability α , draw $i \sim c$ and output (i, i) . Otherwise draw $i \sim q'$ and $\tilde{i} \sim \tilde{q}'$ and output (i, \tilde{i}) . The resulting pair follows a maximal coupling of q and \tilde{q} , since $\mathbb{P}(i = \tilde{i}) = \alpha = 1 - d_{\text{TV}}(q, \tilde{q})$, and marginally $\mathbb{P}(i = n) = \alpha c_n + (1 - \alpha)q'_n = q_n$, and likewise for $\mathbb{P}(\tilde{i} = n)$, for all $n \in \{1, \dots, N\}$. The procedure involves $\mathcal{O}(N)$ operations for N the size of the state space.

We now consider an experiment like those of Yang *et al.* (2016). We define

$$\beta^* = \text{SNR} \sqrt{\left\{ \sigma_0^2 \frac{\log(p)}{n} \right\}} (2, -3, 2, 2, -3, 3, -2, 3, -2, 3, 0, \dots, 0)' \in \mathbb{R}^p,$$

and generate Y given X and β^* from the model with $\sigma^2 = 1$, $\sigma_0^2 = 1$, $n \in \{500, 1000\}$, $p \in \{1000, 5000\}$ and signal-to-noise parameter $\text{SNR} \in \{0.5, 1, 2\}$. We also set $s_0 = 100$, $g = p^3$ and $\kappa = 2$ (exactly as in Yang *et al.* (2016); the value of κ was obtained by personal communication) and generate the covariates X by using a multivariate normal distribution with covariance matrix Σ either equal to a unit diagonal matrix or with entries $\Sigma_{ij} = \exp(-|i - j|)$. We refer to these two cases as the independent design and correlated design cases respectively. We draw from the initial distribution π_0 by creating a vector of p 0s, sampling s_0 co-ordinates uniformly from $\{1, \dots, p\}$ without replacement, and setting the corresponding entries to 1 with probability 0.5.

For various values of n , p and SNR, and the two types of design, we run coupled chains 100 times independently until they meet. We report the average meeting times in Tables 4 and 5. The average meeting times are of the order of 10^4 – 10^5 , depending on the problem; the maximum is attained in the correlated design at $n = 500$, $p = 1000$ and $\text{SNR} = 2$. In contrast with this, the experiments in Yang *et al.* (2016) identify the scenario $n = 500$, $p = 5000$ and $\text{SNR} = 1$ as the most challenging. This discrepancy deserves further study; it could be due to variations from a synthetic data set to another, or to differences in the criteria being reported.

To illustrate the effect of dimension, we focus on the independent design setting with $n = 500$ and $\text{SNR} = 1$, and we consider values of p between 100 and 1000. For each value of p , we run

Table 4. Average meeting times in the independent design

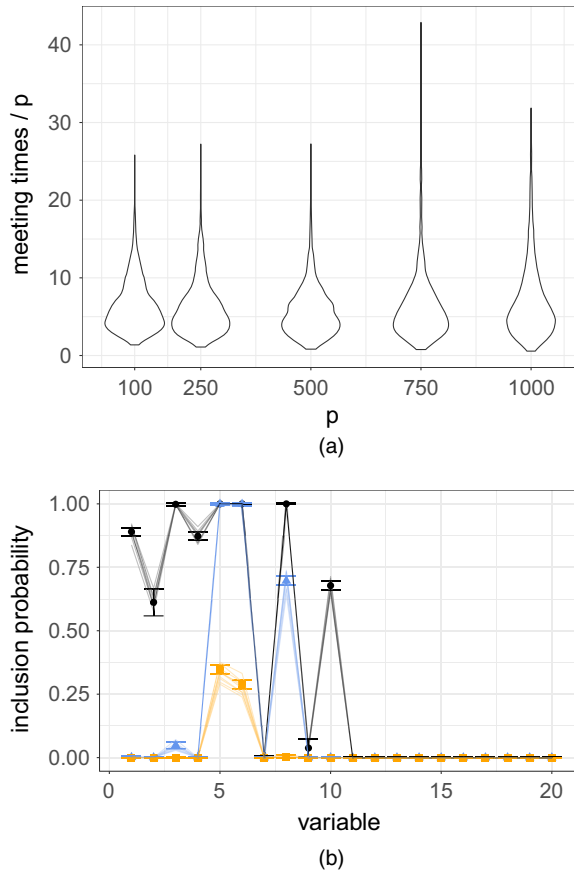
n	p	Results for $\text{SNR} = 0.5$	Results for $\text{SNR} = 1$	Results for $\text{SNR} = 2$
500	1000	4937	7586	6031
500	5000	24634	25602	38083
1000	1000	4729	5893	4892
1000	5000	23407	46398	24712

Table 5. Average meeting times in the correlated design

n	p	Results for $\text{SNR} = 0.5$	Results for $\text{SNR} = 1$	Results for $\text{SNR} = 2$
500	1000	5536	5485	216996
500	5000	27535	28756	29083
1000	1000	4921	5451	5613
1000	5000	24101	29215	23043

1 coupled chains 1000 times independently until they meet. We present violin plots representing
 2 the distributions of meeting times divided by p in Fig. 6(a). The distribution of scaled meeting
 3 times appears to be approximately constant as a function of p , suggesting that meeting times
 4 increase linearly in p . This is consistent with the findings of Yang *et al.* (2016), where mixing
 5 times are shown to increase linearly in p .

6 Focusing now on the independent design case with $n = 500$, $p = 1000$ and $\text{SNR} = 1$ we consider
 7 various values of the prior hyperparameter κ in $\{0.1, 1, 2\}$. We set $k = 75000$ and $m = 150000$
 8 and generate unbiased estimators on a cluster for 120 min, using 200 processors for each value
 9 of κ , and so 600 processors in total. The test function is chosen so that the estimand $\pi(h)$ is
 10 the vector of inclusion probabilities $\mathbb{P}(\gamma_i = 1 | X, Y)$ for $i \in \{1, \dots, 20\}$. Within the time budget,
 11 39282 estimates were produced, with each processor producing between eight and 181 of these.
 12 The largest observed meeting time was 81423. The meeting times were similar across the three
 13 values of κ .



44 **Fig. 6.** (a) Meeting times divided p for $p \in \{100, 250, 500, 750, 1000\}$ and $n = 500$, $\text{SNR} = 1$, in the vari-
 45 able-selection example of Section 5.4 with independent design, based on $R = 1000$ independent repeats (the
 46 violins represent the distributions of scaled meeting times for different p) and (b) posterior probabilities of
 47 inclusion for the first 20 variables, in the setting $n = 500$, $p = 1000$, $\text{SNR} = 1$, for three values of the prior
 48 hyperparameter κ (the error bars representing 95% confidence intervals were obtained after 120 min of
 calculation on 600 processors, using $k = 75000$ and $m = 150000$) (—, —, —, standard MCMC
 estimates based on 10 independent chains of length 10^6 ; ●, $\kappa = 0.1$; ▲, $\kappa = 1$; ■, $\kappa = 2$)

Fig. 6(b) shows the results in the form of 95% confidence intervals shown as error bars, using expression (3.4), the central limit theorem relevant when the time budget is fixed and the number of processors grows large. We observe that κ has a strong effect on the probability of including the first 10 variables in this setting, and that the most satisfactory results are obtained for $\kappa = 0.1$ rather than for $\kappa = 2$, recalling that β^* has non-zero entries in its first 10 components. Note that the error bars are narrow but still noticeable, particularly for $\kappa = 0.1$. On Fig. 6(b), the full lines represent estimates that were obtained with 10 independent MCMC runs with 10^6 iterations each, discarding the first 10^5 iterations as burn-in. These MCMC estimates present noticeable variability in spite of the large number of iterations. In a standard MCMC setting, we might run chains for more iterations until the estimates agree across independent runs. In the framework proposed, we increase the precision by generating more independent unbiased estimators without necessarily modifying k or m .

Fig. 6(b) suggests that the variable selection procedure that is considered here is sensitive to the prior hyperparameter κ ; we refer to Yang *et al.* (2016), and to Johnson (2013) and Nikooienjad *et al.* (2016) for related discussions on Bayesian variable selection in high dimension and convergence of MCMC algorithms.

5.5. Cut distribution

Finally, our proposed estimator can be used to approximate the cut distribution, which poses a significant challenge for existing MCMC methods (Plummer, 2014; Jacob *et al.*, 2017). This illustrates another appeal of the unbiasedness property, beyond the motivation for parallel computation.

Consider two models: one with parameters θ_1 and data Y_1 and another with parameters θ_2 and data Y_2 , where the likelihood of Y_2 might depend on both θ_1 and θ_2 . For instance the first model could be a regression with data Y_1 and coefficients θ_1 , and the second model could be another regression whose covariates are the residuals, coefficients or fitted values of the first regression (Pagan, 1984; Murphy and Topel, 2002). In principle one could introduce an encompassing model and conduct joint inference on θ_1 and θ_2 via the posterior distribution. In that case, misspecification of either model would lead to misspecification of the ensemble and thus to a misleading quantification of uncertainty, as noted in several studies (e.g. Liu *et al.* (2009), Plummer (2014), Lunn *et al.* (2009), McCandless *et al.* (2010), Zigler (2016) and Blangiardo *et al.* (2011)).

The cut distribution (Spiegelhalter *et al.*, 2003; Plummer, 2014) allows the propagation of uncertainty about θ_1 to inference on θ_2 while preventing misspecification in the second model from affecting estimation in the first. The cut distribution is defined as

$$\pi_{\text{cut}}(\theta_1, \theta_2) = \pi_1(\theta_1)\pi_2(\theta_2|\theta_1).$$

Here $\pi_1(\theta_1)$ refers to the distribution of θ_1 given Y_1 in the first model alone, and $\pi_2(\theta_2|\theta_1)$ refers to the distribution of θ_2 given Y_2 and θ_1 in the second model. Often, the density $\pi_2(\theta_2|\theta_1)$ can be evaluated up to only a constant in θ_2 , which may vary with θ_1 . This makes the cut distribution difficult to approximate with MCMC algorithms (Plummer, 2014).

A naive approach consists of first running an MCMC algorithm targeting $\pi_1(\theta_1)$ to obtain a sample $(\theta_1^n)_{n=1}^{N_1}$, perhaps after discarding a burn-in period and thinning the chain. Then, for each θ_1^n , one can run an MCMC algorithm targeting $\pi_2(\theta_2|\theta_1^n)$, yielding N_2 samples. One might again discard some burn-in and thin the chains, or just keep the final state of each chain. The resulting joint samples approximate the cut distribution. However, the validity of this approach

relies on a double limit in N_1 and N_2 . Diagnosing convergence may also be difficult given the number of chains in the second stage, each of which targets a different distribution $\pi_2(\theta_2|\theta_1^n)$.

If we could sample $\theta_1 \sim \pi_1$ and $\theta_2|\theta_1 \sim \pi_2(\theta_2|\theta_1)$, then the pair (θ_1, θ_2) would follow the cut distribution. The same two-stage rationale can be applied in the framework proposed. Consider a test function $(\theta_1, \theta_2) \mapsto h(\theta_1, \theta_2)$. Writing \mathbb{E}_{cut} for expectations with respect to π_{cut} , the law of iterated expectations yields

$$\mathbb{E}_{\text{cut}}[h(\theta_1, \theta_2)] = \int \left\{ \int h(\theta_1, \theta_2) \pi_2(d\theta_2|\theta_1) \right\} \pi_1(d\theta_1) = \int \bar{h}(\theta_1) \pi_1(d\theta_1).$$

Here $\bar{h}(\theta_1) = \int h(\theta_1, \theta_2) \pi_2(d\theta_2|\theta_1)$. In the framework proposed, we can make an unbiased estimator of $\bar{h}(\theta_1)$ for all θ_1 and then plug these estimators into an unbiased estimator of the integral $\int \bar{h}(\theta_1) \pi_1(d\theta_1)$. This is perhaps clearer by using the signed measure representation of Section 2.4: one can obtain a signed measure $\hat{\pi}_1 = \sum_{l=1}^N \omega_l \delta_{\theta_{1,l}}$ approximating π_1 , and then obtain an unbiased estimator of $\bar{h}(\theta_{1,l})$ for all l , denoted by \bar{H} . Then the weighted average $\sum_{l=1}^N \omega_l \bar{H}$ is an unbiased estimator of $\mathbb{E}_{\text{cut}}[h(\theta_1, \theta_2)]$ by the law of iterated expectations. Such estimators can be generated independently in parallel, and their average provides a consistent approximation of an expectation with respect to the cut distribution.

We consider the example that was described in Plummer (2014), inspired by an investigation of the international correlation between human papilloma virus (HPV) prevalence and cervical cancer incidence (Maucourt-Boulch *et al.*, 2008). The first module concerns HPV prevalence, with data independently collected in 13 countries. The parameter $\theta_1 = (\theta_{1,1}, \dots, \theta_{1,13})$ receives a beta(1,1) prior distribution independently for each component. The data (Y_1, \dots, Y_{13}) consist of 13 pairs of integers. The first represents the number of women who were infected with high-risk HPV, and the second represents population sizes. The likelihood specifies a binomial model for Y_i , independently for each component i . The posterior for this model is given by a product of beta distributions.

The second module concerns the relationship between HPV prevalence and cancer incidence, and posits a Poisson regression. The parameters are $\theta_2 = (\theta_{2,1}, \theta_{2,2}) \in \mathbb{R}^2$ and receive a normal prior with zero mean and variance 10^3 per component. The likelihood in this module is given by

$$Z_{1,i} \sim \text{Poisson}\{\exp(\theta_{2,1} + \theta_{1,i}\theta_{2,2} + Z_{2,i})\} \quad \text{for } i \in \{1, \dots, 13\},$$

where the data $(Z_{1,i}, Z_{2,i})_{i=1}^{13}$ are pairs of integers. The first component represents numbers of cancer cases, whereas the second is the number of woman-years of follow-up. The Poisson regression model might be misspecified, motivating departures from inference based on the joint model (Plummer, 2014).

Here we can draw directly from the first posterior, denoted by $\pi_1(\theta_1)$, and obtain a sample $(\theta_1^n)_{n=1}^N$. For each θ_1^n we consider an MH algorithm targeting $\pi_2(\theta_2|\theta_1^n)$, using a normal random-walk proposal with variance Σ . We couple this algorithm by using reflection maximal couplings of the proposals as in Section 4.1. In preliminary runs, starting with a standard bivariate normal distribution as an initial distribution and a proposal covariance matrix set to identity, we estimate the first two moments of the cut distribution, and we use them to refine the initial distribution π_0 and the proposal covariance matrix Σ . With these settings we obtain the distribution of meeting times that is shown in Fig. 7(a). We then set $k = 100$ and $m = 10k$, and obtain approximations of the cut distribution represented by histograms in Figs 7(b) and 7(c), using $N = 10000$ unbiased estimators. The overlaid curves correspond to a kernel density estimate obtained by running $m = 1000$ steps of MCMC targeting $\pi_2(\theta_2|\theta_1^n)$ with θ_1^n drawn from $\pi_1(\theta_1)$, for $n \in \{1, \dots, N\}$, and keeping the final m th state of each chain. The proposed estimators can be refined by increasing

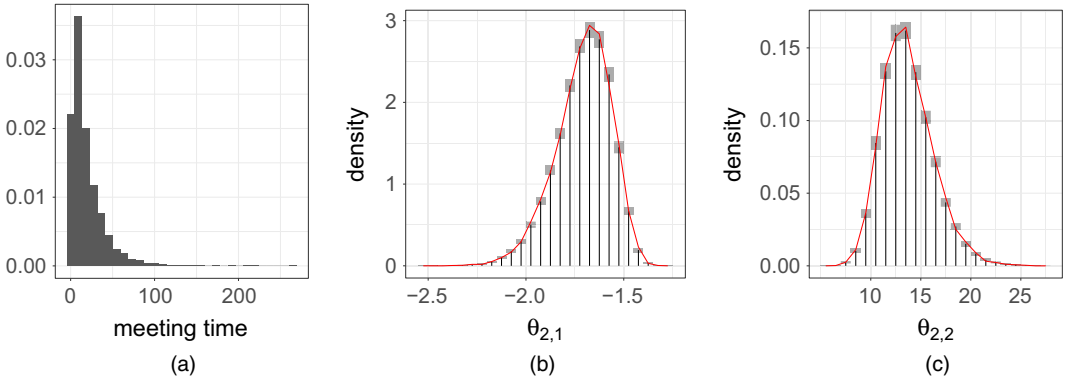


Fig. 7. (a) Meeting times and histograms of (b) $\theta_{2,1}$ and (c) $\theta_{2,2}$, in the example of Section 5.5, computed from 10000 unbiased estimators, with $k = 100$ and $m = 1000$ (■, 95% confidence intervals on the histogram estimates; —, marginal probability density functions of the cut distribution, obtained by running $N = 10000$ MCMC chains for 1000 steps, each targeting $\pi_2(\theta_2|\theta_1^q)$ for θ_1^q drawn from $\pi_1(\theta_1)$, and retaining the last state only)

the number N of independent replications, whereas the MCMC estimators would converge only in the double limit of N and m going to ∞ .

6. Discussion

By combining the powerful technique of Glynn and Rhee (2014) with couplings of MCMC algorithms, unbiased estimators of integrals with respect to the target distribution can be constructed. Their efficiency can be controlled with tuning parameters k and m , for which we have proposed guidelines: k can be chosen as a large quantile of the meeting time τ , and m as a multiple of k . Improving on these simple guidelines stands as a subject for future research. In numerical experiments we have argued that the estimators proposed yield a practical way of parallelizing MCMC computations in a range of settings. We stress that coupling pairs of Markov chains does not improve their marginal mixing properties, and that poor mixing of the underlying chains can lead to poor performance of the resulting estimator. The choice of initial distribution π_0 can have undesirable effects on the estimators, as in the multimodal example of Section 5.1. Unreliable estimators would also result from stopping the chains before their meeting time.

Couplings of MCMC algorithms can be devised by using maximal couplings reflection couplings and common random numbers. We have focused on couplings that can be implemented without further analytical knowledge about the target distribution or about the MCMC kernels. However, these couplings might result in prohibitively large meeting times, either because the marginal chains mix slowly, as in Section 5.1, or because the coupling strategy is ineffective, as in Section 4.2.

Regarding convergence diagnostics, the framework proposed yields the following representation for the total variation between π_k and π , where π_k denotes the marginal distribution of X_k :

$$\begin{aligned} d_{\text{TV}}(\pi_k, \pi) &= \frac{1}{2} \sup_{h: |h| \leq 1} |\mathbb{E}[h(X_k)] - \mathbb{E}_\pi[h(X)]| \\ &= \frac{1}{2} \sup_{h: |h| \leq 1} \left| \mathbb{E} \left[\sum_{t=k+1}^{\tau-1} \{h(X_t) - h(Y_{t-1})\} \right] \right|, \end{aligned}$$

Here the supremum ranges over all bounded measurable functions under the assumptions stated. The equality above has several consequences. For instance, the triangle inequality implies that $d_{TV}(\pi_k, \pi) \leq \min(1, \mathbb{E}[\max\{0, (\tau - k - 1)\}])$, and we can approximate $\mathbb{E}[\max\{0, (\tau - k - 1)\}]$ by Monte Carlo sampling for a range of k -values. This is pursued in Biswas and Jacob (2019), where the construction proposed is extended to allow for arbitrary time lags between the coupled chains.

Thanks to its potential for parallelization, the framework proposed can facilitate a consideration of MCMC kernels that might be too expensive for serial implementation. For instance, one can improve MH-within-Gibbs samplers by performing more MH steps per component update, Hamiltonian Monte Carlo sampling by using smaller step sizes in the numerical integrator (Heng and Jacob, 2019) and particle MCMC sampling by using more particles in the particle filters (Andrieu *et al.*, 2010; Jacob *et al.*, 2019). We expect the optimal tuning of MCMC kernels to be different in the proposed framework from when used marginally.

On top of enabling the application of the results of Glynn and Heidelberger (1991) to accommodate budget constraints, the lack of bias of the estimators proposed can be beneficial in combination with the law of total expectation, to implement modular inference procedures as in Section 5.5. In Rischard *et al.* (2018) the lack of bias was exploited in new estimators of Bayesian cross-validation criteria. In Chen *et al.* (2018) similar unbiased estimators were used in the expectation step of an expectation–maximization algorithm. There may be other settings where the lack of bias is appealing, for instance in gradient estimation for stochastic gradient descents (Tadić *et al.*, 2017).

Acknowledgements

The authors are grateful to Jeremy Heng and Luc Vincent-Genod for useful discussions. The authors gratefully acknowledge support by the National Science Foundation through grants DMS-1712872 (Pierre E. Jacob) and DMS-1513040 (Yves F. Atchadé).

References

- Altekar, G., Dwarkadas, S., Huelsenbeck, J. P. and Ronquist, F. (2004) Parallel Metropolis coupled Markov chain Monte Carlo for Bayesian phylogenetic inference. *Bioinformatics*, **20**, 407–415.
- Andrieu, C., Doucet, A. and Holenstein, R. (2010) Particle Markov chain Monte Carlo methods (with discussion). *J. R. Statist. Soc. B*, **72**, 269–342.
- Andrieu, C. and Vihola, M. (2015) Convergence properties of pseudo-marginal Markov chain Monte Carlo algorithms. *Ann. Appl. Probab.*, **25**, 1030–1077.
- Atchadé, Y. F. (2006) An adaptive version for the Metropolis adjusted Langevin algorithm with a truncated drift. *Methodol. Comput. Appl. Probab.*, **8**, 235–254.
- Atchadé, Y. F. (2016) Markov chain Monte Carlo confidence intervals. *Bernoulli*, **22**, 1808–1838.
- Biswas, N. and Jacob, P. E. (2019) Estimating convergence of Markov chains with L-lag couplings. *Preprint arXiv:1905.09971*.
- Blangiardo, M., Hansell, A. and Richardson, S. (2011) A Bayesian model of time activity data to investigate health effect of air pollution in time series studies. *Atmosph. Environ.*, **45**, 379–386.
- Bou-Rabee, N., Eberle, A. and Zimmer, R. (2018) Coupling and convergence for Hamiltonian Monte Carlo. *Preprint arXiv:1805.00452*.
- Brockwell, A. E. (2006) Parallel Markov chain Monte Carlo simulation by pre-fetching. *J. Computn. Graph. Statist.*, **15**, 246–261.
- Brockwell, A. E. and Kadane, J. B. (2005) Identification of regeneration times in MCMC simulation, with application to adaptive schemes. *J. Computn. Graph. Statist.*, **14**, 436–458.
- Brooks, S. P., Gelman, A., Jones, G. and Meng, X.-L. (2011) *Handbook of Markov chain Monte Carlo*. Boca Raton: CRC Press.
- Calderhead, B. (2014) A general construction for parallelizing Metropolis–Hastings algorithms. *Proc. Natn. Acad. Sci.*, **111**, 17408–17413.
- Casella, G., Lavine, M. and Robert, C. P. (2001) Explaining the perfect sampler. *Am. Statistn.*, **55**, 299–305.

- Chatterjee, S., Guntuboyina, A., Sen, B. *et al.* (2015) On risk bounds in isotonic and other shape restricted regression problems. *Ann. Statist.*, **43**, 1774–1800.
- Chen, W., Ma, L. and Liang, X. (2018) Blind identification based on expectation-maximization algorithm coupled with blocked Rhee–Glynn smoothing estimator. *IEEE Commun. Lett.*, **22**, 1838–1841.
- Choi, H. M. and Hobert, J. P. (2013) The Pólya–Gamma Gibbs sampler for Bayesian logistic regression is uniformly ergodic. *Electron. J. Statist.*, **7**, 2054–2064.
- Cowles, M. K. and Rosenthal, J. S. (1998) A simulation approach to convergence rates for Markov chain Monte Carlo algorithms. *Statist. Comput.*, **8**, 115–124.
- Diaconis, P. and Freedman, D. (1999) Iterated random functions. *SIAM Rev.*, **41**, 45–76.
- Douc, R., Moulines, E. and Rosenthal, J. S. (2004) Quantitative bounds on convergence of time-inhomogeneous Markov chains. *Ann. Appl. Probab.*, **14**, 1643–1665.
- Duane, S., Kennedy, A. D., Pendleton, B. J. and Roweth, D. (1987) Hybrid Monte Carlo. *Phys. Lett. B*, **195**, 216–222.
- Flegal, J. M., Haran, M. and Jones, G. L. (2008) Markov chain Monte Carlo: can we trust the third significant figure? *Statist. Sci.*, 250–260.
- Flegal, J. M. and Herbei, R. (2012) Exact sampling for intractable probability distributions via a Bernoulli factory. *Electron. J. Statist.*, **6**, 10–37.
- Gaver, D. P. and O’Muircheartaigh, I. G. (1987) Robust empirical Bayes analyses of event rates. *Technometrics*, **29**, 1–15.
- Geyer, C. J. (1991) Markov chain Monte Carlo maximum likelihood. *Technical Report*. School of Statistics, University of Minnesota.
- Glynn, P. W. (2016) Exact simulation versus exact estimation. In *Proc. Winter Simulation Conf*, pp. 193–205. New York: Institute of Electrical and Electronics Engineers.
- Glynn, P. W. and Heidelberger, P. (1990) Bias properties of budget constraint simulations. *Ops Res.*, **38**, 801–814.
- Glynn, P. W. and Heidelberger, P. (1991) Analysis of parallel replicated simulations under a completion time constraint. *ACM Trans. Model. Comput. Simulns*, **1**, 3–23.
- Glynn, P. W. and Rhee, C.-H. (2014) Exact estimation for Markov chain equilibrium expectations. *J. Appl. Probab.*, **51**, A, 377–389.
- Glynn, P. W. and Whitt, W. (1992) The asymptotic efficiency of simulation estimators. *Ops Res.*, **40**, 505–520.
- Gong, L. and Flegal, J. M. (2016) A practical sequential stopping rule for high-dimensional Markov chain Monte Carlo. *J. Computat. Graph. Statist.*, **25**, 684–700.
- Goodman, J., Weare, J. *et al.* (2010) Ensemble samplers with affine invariance. *Commun. Appl. Math. Computat. Sci.*, **5**, 65–80.
- Goudie, R. J., Turner, R. M., De Angelis, D. and Thomas, A. (2017) Massively parallel MCMC for Bayesian hierarchical models. *Preprint arXiv:1704.03216*.
- Green, P. J., Łatuszyński, K., Pereyra, M. and Robert, C. P. (2015) Bayesian computation: a summary of the current state, and samples backwards and forwards. *Statist. Comput.*, **25**, 835–862.
- Hastings, W. K. (1970) Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, **57**, 97–109.
- Heng, J. and Jacob, P. E. (2019) Unbiased Hamiltonian Monte Carlo with couplings. *Biometrika*, **106**, 287–302.
- Hoffman, M. D. and Gelman, A. (2014) The No-U-Turn Sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *J. Mach. Learn. Res.*, **15**, 1593–1623.
- Huber, M. (2016) *Perfect Simulation*. Boca Raton: CRC Press.
- Jacob, P. E., Lindsten, F. and Schön, T. B. (2019) Smoothing with couplings of conditional particle filters. *J. Am. Statist. Ass.*, to be published, doi 10.1080/01621459.2018.1548856.
- Jacob, P. E., Murray, L. M., Holmes, C. C. and Robert, C. P. (2017) Better together?: Statistical learning in models made of modules. *Preprint arXiv:1708.08719*.
- Jacob, P. E., Robert, C. P. and Smith, M. H. (2011) Using parallel computation to improve independent Metropolis–Hastings based estimation. *J. Computat. Graph. Statist.*, **20**, 616–635.
- Jarner, S. F. and Hansen, E. (2000) Geometric ergodicity of Metropolis algorithms. *Stoch. Processes Appl.*, **85**, 341–361.
- Johndrow, J. E. and Mattingly, J. C. (2017) Coupling and decoupling to bound an approximating Markov chain. *Preprint arXiv:1706.02040*.
- Johnson, V. E. (1996) Studying convergence of Markov chain Monte Carlo algorithms using coupled sample paths. *J. Am. Statist. Ass.*, **91**, 154–166.
- Johnson, V. E. (1998) A coupling-regeneration scheme for diagnosing convergence in Markov chain Monte Carlo algorithms. *J. Am. Statist. Ass.*, **93**, 238–248.
- Johnson, V. E. (2013) On numerical aspects of Bayesian model selection in high and ultrahigh-dimensional settings. *Baysn Anal.*, **8**, 741–758.
- Khare, K. and Hobert, J. P. (2013) Geometric ergodicity of the Bayesian lasso. *Electron. J. Statist.*, **7**, 2150–2163.
- Lee, A., Doucet, A. and Łatuszyński, K. (2014) Perfect simulation using atomic regeneration with application to sequential Monte Carlo. *Preprint arXiv:1407.5770*.

- 1 Lee, A., Yau, C., Giles, M. B., Doucet, A. and Holmes, C. C. (2010) On the utility of graphics cards to perform
 2 massively parallel simulation of advanced Monte Carlo methods. *J. Computnl Graph. Statist.*, **19**, 769–789.
- 3 Lindvall, T. (2002) *Lectures on the Coupling Method*. Dover Books.
- 4 Liu, J. S. (2008) *Monte Carlo Strategies in Scientific Computing*. New York: Springer Science and Business Media.
- 5 Liu, F., Bayarri, M. and Berger, J. (2009) Modularization in Bayesian analysis, with emphasis on analysis of
 6 computer models. *Bayesn Anal.*, **4**, 119–150.
- 7 Lunn, D., Best, N., Spiegelhalter, D., Graham, G. and Neuwenschwander, B. (2009) Combining MCMC with
 8 sequential PKPD modelling. *J. Pharmkinet. Pharmdynam.*, **36**, 19–38.
- 9 Mainini, E. (2012) A description of transport cost for signed measures. *J. Math. Sci.*, **181**, 837–855.
- 10 Mangoubi, O. and Smith, A. (2017) Rapid mixing of Hamiltonian Monte Carlo on strongly log-concave distri-
 11 butions. *Preprint arXiv:1708.07114*.
- 12 Maucort-Boulch, D., Franceschi, S. and Plummer, M. (2008) International correlation between human papillo-
 13 mavirus prevalence and cervical cancer incidence. *Cancer Epidem. Biomark. Prevn.*, **17**, 717–720.
- 14 McCandless, L. C., Douglas, I. J., Evans, S. J. and Smeeth, L. (2010) Cutting feedback in Bayesian regression
 15 adjustment for the propensity score. *Int. J. Biostatist.*, **6**, no. 2.
- 16 McLeish, D. (2011) A general method for debiasing a Monte Carlo estimator. *Monte Carlo Meth. Appl.*, **17**,
 17 301–315.
- 18 Meyn, S. and Tweedie, R. (2009) *Markov Chains and Stochastic Stability*, 2nd edn. Cambridge: Cambridge Uni-
 19 versity Press.
- 20 Middleton, L., Deligiannidis, G., Doucet, A. and Jacob, P. E. (2018) Unbiased Markov chain Monte Carlo for
 21 intractable target distributions. *Preprint arXiv:1807.08691*.
- 22 Middleton, L., Deligiannidis, G., Doucet, A. and Jacob, P. E. (2019) Unbiased smoothing using particle independ-
 23 ent Metropolis-Hastings. *Proc. Mach. Learn. Res.*, **89**, 2378–2387.
- 24 Murdoch, D. J. and Green, P. J. (1998) Exact sampling from a continuous state space. *Scand. J. Statist.*, **25**,
 25 483–502.
- 26 Murphy, K. M. and Topel, R. H. (2002) Estimation and inference in two-step econometric models. *J. Bus. Econ.*
 27 *Statist.*, **20**, 88–97.
- 28 Murray, I., Adams, R. P. and MacKay, D. J. C. (2010) Elliptical slice sampling. In *Proc. 13th Int. Conf. Artificial*
 29 *Intelligence and Statistics*, pp. 541–548.
- 30 Mykland, P., Tierney, L. and Yu, B. (1995) Regeneration in Markov chain samplers. *J. Am. Statist. Ass.*, **90**,
 31 233–241.
- 32 Neal, R. M. (1993) Bayesian learning via stochastic dynamics. In *Advances in Neural Information Processing*
 33 *Systems*, pp. 475–475.
- 34 Neal, R. M. (1999) Circularly-coupled Markov chain sampling. In *Technical Report*. Department of Statistics,
 35 University of Toronto, Toronto.
- 36 Neal, R. M. (2011) MCMC using Hamiltonian dynamics. In *Handbook of Markov chain Monte Carlo*, vol. **2**.
- 37 Neal, R. M. and Pinto, R. L. (2001) Improving Markov chain Monte Carlo estimators by coupling to an approx-
 38 imating chain. *Technical Report*. Department of Statistics, University of Toronto, Toronto.
- 39 Nikooienjad, A., Wang, W. and Johnson, V. E. (2016) Bayesian variable selection for binary outcomes in high-
 40 dimensional genomic studies using non-local priors. *Bioinformatics*, **32**, 1338–1345.
- 41 Nummelin, E. (2002) MC’s for MCMC’ists. *Int. Statist. Rev.*, **70**, 215–240.
- 42 Owen, A. B. (2017) Statistically efficient thinning of a Markov chain sampler. *J. Computnl Graph. Statist.*, **26**,
 43 738–744.
- 44 Pagan, A. (1984) Econometric issues in the analysis of regressions with generated regressors. *Int. Econ. Rev.*,
 45 221–247.
- 46 Pal, S. and Khare, K. (2014) Geometric ergodicity for Bayesian shrinkage models. *Electron. J. Statist.*, **8**, 604–645.
- 47 Plummer, M. (2014) Cuts in Bayesian graphical models. *Statist. Comput.*, 1–7.
- 48 Plummer, M., Best, N., Cowles, K. and Vines, K. (2006) CODA: convergence diagnosis and output analysis for
 MCMC. *R News*, **6**, 7–11.
- Plummer, M. *et al.* (2003) JAGS: a program for analysis of Bayesian graphical models using Gibbs sampling. In
Proc. 3rd Int. Wrkshp Distributed Statistical Computin, Vienna.
- Pollock, M., Fearnhead, P., Johansen, A. M. and Roberts, G. O. (2016) The scalable Langevin exact algorithm:
 Bayesian inference for big data. *Preprint arXiv:1609.03436*.
- Propp, J. G. and Wilson, D. B. (1996) Exact sampling with coupled Markov chains and applications to statistical
 mechanics. *Rand. Struct. Algrthms*, **9**, 223–252.
- R Core Team (2015) *R: a Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical
 Computing.
- Reutter, A. and Johnson, V. E. (1995) General strategies for assessing convergence of MCMC algorithms using
 coupled sample paths. Institute of Statistics and Decision Sciences, Duke University, Durham.
- Rhee, C.-H. and Glynn, P. W. (2012) A new approach to unbiased estimation for SDE’s. In *Proc. Winter Simulation*
Conf., 17.
- Rischar, M., Jacob, P. E. and Pillai, N. (2018) Unbiased estimation of log normalizing constants with applications
 to Bayesian cross-validation. *Preprint arXiv:1810.01382*.

- 1 Robert, C. P. and Casella, G. (2004) *Monte Carlo Statistical Methods*, 2nd edn. New York: Springer.
- 2 Roberts, G. O., Gelman, A. and Gilks, W. R. (1997) Weak convergence and optimal scaling of random walk
3 Metropolis algorithms. *Ann. Appl. Probab.*, **7**, 110–120.
- 4 Roberts, G. O. and Rosenthal, J. S. (2004) General state space Markov chains and MCMC algorithms. *Probab.*
5 *Surv.*, **1**, 20–71.
- 6 Roberts, G. O. and Tweedie, R. L. (1996a) Exponential convergence of Langevin distributions and their discrete
7 approximations. *Bernoulli*, 341–363.
- 8 Roberts, G. O. and Tweedie, R. L. (1996b) Geometric convergence and central limit theorems for multidimensional
9 Hastings and Metropolis algorithms. *Biometrika*, **83**, 95–110.
- 10 Rosenthal, J. S. (1996) Analysis of the Gibbs sampler for a model related to James-Stein estimators. *Statist.*
11 *Comput.*, **6**, 269–275.
- 12 Rosenthal, J. S. (1997) Faithful couplings of Markov chains: now equals forever. *Adv. Appl. Math.*, **18**, 372–381.
- 13 Rosenthal, J. S. (2000) Parallel computing and Monte Carlo algorithms. *Far East J. Theoret. Statist.*, **4**, 207–236.
- 14 Rosenthal, J. S. (2002) Quantitative convergence rates of Markov chains: a simple account. *Electron. Commun.*
15 *Probab.*, **7**, 123–128.
- 16 Spiegelhalter, D., Thomas, A., Best, N. and Lunn, D. (2003) *WinBUGS User Manual*.
- 17 Srivastava, S., Cevher, V., Dinh, Q. and Dunson, D. (2015) WASP: scalable Bayes via barycenters of subset
18 posteriors. In *Artificial Intelligence and Statistics*, pp. 912–920.
- 19 Swendsen, R. H. and Wang, J.-S. (1987) Nonuniversal critical dynamics in Monte Carlo simulations. *Phys. Rev.*
20 *Lett.*, **58**, 86.
- 21 Tadić, V. B., Doucet, A. *et al.* (2017) Asymptotic bias of stochastic gradient search. *Ann. Appl. Probab.*, **27**,
22 3255–3304.
- 23 Thorisson, H. (2000) *Coupling, Stationarity, and Regeneration*. New York: Springer.
- 24 Titsias, M. K. and Yau, C. (2017) The Hamming ball sampler. *J. Am. Statist. Ass.*, 1–14.
- 25 Tjelmeland, H. (2004) Using all Metropolis–Hastings proposals to estimate mean values. *Technical Report*. De-
26 partment of Mathematical Sciences, Norwegian University of Science and Technology.
- 27 Tweedie, R. (1983) The existence of moments for stationary Markov chains. *J. Appl. Probab.*, **20**, 191–196.
- 28 Vanetti, P., Bouchard-Côté, A., Deligiannidis, G. and Doucet, A. (2017) Piecewise deterministic Markov chain
29 Monte Carlo. *Preprint arXiv:1707.05296*.
- 30 Vats, D., Flegal, J. M., Jones, G. L. *et al.* (2018) Strong consistency of multivariate spectral variance estimators
31 in Markov chain Monte Carlo. *Bernoulli*, **24**, 1860–1909.
- 32 Wang, X., Guo, F., Heller, K. A. and Dunson, D. B. (2015) Parallelizing MCMC with random partition trees. In
33 *Advances in Neural Information Processing Systems*, pp. 451–459.
- 34 Wolff, U. (1989) Comparison between cluster Monte Carlo algorithms in the Ising model. *Phys. Lett. B*, **228**,
35 379–382.
- 36 Yang, S., Chen, Y., Bernton, E. and Liu, J. S. (2017) On parallelizable Markov chain Monte Carlo algorithms
37 with waste-recycling. *Statist. Comput.*, 1–9.
- 38 Yang, Y., Wainwright, M. J. and Jordan, M. I. (2016) On the computational complexity of high-dimensional
39 Bayesian variable selection. *Ann. Statist.*, **44**, 2497–2532.
- 40 Zigler, C. M. (2016) The central role of Bayes theorem for joint estimation of causal effects and propensity scores.
41 *Am. Statistn.*, **70**, 47–54.

Supporting information

42 Additional ‘supporting information’ may be found in the on-line version of this article:

43 ‘Unbiased Markov chain Monte Carlo with couplings’.