

Unbiased Estimation using Underdamped Langevin Dynamics

HAMZA RUZAYQAT, NEIL K. CHADA, & AJAY JASRA

Applied Mathematics and Computational Science Program,
Computer, Electrical and Mathematical Sciences and Engineering Division,

King Abdullah University of Science and Technology, Thuwal, 23955-6900, KSA.

E-Mail: hamza.ruzayqat@kaust.edu.sa, neilchada123@gmail.com, ajay.jasra@kaust.edu.sa

Abstract

In this work we consider the unbiased estimation of expectations w.r.t. probability measures that have non-negative Lebesgue density, and which are known point-wise up-to a normalizing constant. We focus upon developing an unbiased method via the underdamped Langevin dynamics, which has proven to be popular of late due to applications in statistics and machine learning. Specifically in continuous-time, the dynamics can be constructed so that as the time goes to infinity they admit the probability of interest as a stationary measure. In many cases, time-discretized versions of the underdamped Langevin dynamics are used in practice which are run only with a fixed number of iterations. We develop a novel scheme based upon doubly randomized estimation as in [17, 19], which requires access only to time-discretized versions of the dynamics. The proposed scheme aims to remove the discretization bias and the bias resulting from running the dynamics for a finite number of iterations. We prove, under standard assumptions, that our estimator is of finite variance and either has finite expected cost, or has finite cost with a high probability. To illustrate our theoretical findings we provide numerical experiments which verify our theory, which include challenging examples from Bayesian statistics and statistical physics.

Key words: Underdamped Langevin dynamics, unbiased estimation, maximal couplings, Markov chain simulation

AMS subject classifications: 60J22, 65C05, 65C40, 82C31, 62G08, 35Q56

Code available at: https://github.com/ruzayqat/unbiased_uld

Corresponding author: Hamza Ruzayqat. E-mail: hamza.ruzayqat@kaust.edu.sa

1 Introduction

We consider a class of probability measures on the measurable space $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ with non-negative Lebesgue densities known point-wise up-to a normalizing constant. The objective of this article is to consider simulation-based methods, which can return stochastic estimates of finite expectations of functions w.r.t. the afore-mentioned probability measure, that are unbiased, that is, on average are equal to the expectation of interest. This latter task is of interest in a wide variety of areas such as applied mathematics, physics and statistics; see e.g. [36] for a book length introduction.

The primary methodology that is used in the literature to approximate expectations is that of Markov chain Monte Carlo (MCMC) methods. These are schemes which generate ergodic Markov chains whose stationary distribution is exactly the one of interest and there are numerous variants of MCMC, such as random walk Metropolis-Hastings, Hamiltonian Monte Carlo and non-reversible MCMC; see e.g. [36]. In addition to this, are methods based upon uncorrected time discretizations of continuous-time processes, which also have the appropriate distribution of interest as a stationary

distribution, such as the underdamped Langevin algorithm e.g. [13, 32]. In particular with the underdamped Langevin algorithm, it has been empirically observed to converge with a better rate, to an invariant distribution, than that of the overdamped Langevin dynamics, which is much simpler in comparison [7, 9, 12]. As a result, these latter methods are of interest often due to their relative ease of simulation relative to the former and have gained significant popularity in the statistics and machine learning literature [14, 41]. In both of the classes of methods that we have mentioned, in general, without starting these Markov chains from draws from the target distribution of interest, one seldom returns unbiased estimates which is the main interest in this article. Unbiasedness can be desirable in certain contexts, for instance when computing sensitivities for stochastic gradient algorithms e.g. [1].

We focus upon contributing to the class of unbiased MCMC algorithms and trying to enhance the applicability of such schemes. Unbiased estimation can be achieved in at least two ways, by exact simulation [2, 3] (e.g. starting the chain from the correct probability of interest) or unbiased approximation. The former has been investigated many years ago in the guise of coupling from the past MCMC algorithms (e.g. [33]), but due to the complications of doing so, the application of such methods is rather rare. The latter is often based upon the pioneering works on unbiased estimation found in [16] (see also [29, 35]). The idea of these latter methods, in the Markov chain context, is to work with a pair of Markov chains on a product space and simulate them until they are equal (the *meeting time*); there is then a novel identity which ensures that the estimate is unbiased. Creating a methodology for allowing the Markov chains to meet was developed in the paper [21] and subsequent to this, several modifications and improvements were given in [30, 43]. However there has been recent interest of unbiased methodologies, related to that latter way, where we provide some of these works [5, 6, 17, 39, 40], which made extensions in the context of particle filtering and particle MCMC methods.

One of the main issues of the methodology developed in [21] and the sequels, is that one needs to consider sometimes quite complex coupling of pairs of Markov chains. This can be quite non-trivial to achieve and, in some scenarios such as when the target density is multimodal, rather inefficient, leading to large variances in estimation. This was partially addressed in [38] which considered an unbiased version of the Schrödinger-Föllmer sampler (SFS). The latter is a diffusion process on a bounded time domain $[0, 1]$, that transports a degenerate distribution at 0, to the target of interest, assuming that the latter is absolutely continuous w.r.t. a d -dimensional standard Gaussian. Even under Euler time-discretization, the process cannot be simulated as the drift term is complicated resulting it in being intractable, but several mechanisms are available; see [38]. The authors in that paper show that by using doubly randomized unbiased schemes (e.g. [22, 23]) an unbiased version of the SFS can be developed which provides unbiased estimates without having to resort to complex coupling mechanisms. One of the drawbacks of the methodology, however, is the SFS method in the beginning; when approximating the drift, the method can struggle to well-represent complex probability measures.

As a result the focus of this article is on the development of unbiased schemes which alleviate the issues discussed above. Specifically we want to consider unbiased schemes, which can handle two forms of bias, that from the MCMC and from the discretization bias, arising from models such as stochastic differential equations. This extends the work of [21] which does not consider the latter as a form of bias. To do so we exploit ideas randomized multilevel Monte Carlo (MLMC) methods, which have been developed to reduce the cost to attain a particular order of MSE $\mathcal{O}(\epsilon^2)$, $\epsilon > 0$. A recent study showing the connection with the unbiased MC and MLMC is the work of Vihola [44].

Also we aim to provide a method which does not require a complex coupling of Markov chains, which is relatively simple to implement, which has been well understood in theory and practice.

1.1 Contributions

The contribution of this article is to develop a new version of the underdamped Langevin algorithm which, even when only working with time-discretizations of the process, can deliver unbiased estimates of expectations w.r.t. the class of probability measures under considerations. Our motivation is that, like similarly to the work in [38], it does not require very complex coupling techniques and also relies on a double randomization scheme which was developed in [19]. The method also retains the advantages of the unbiased MCMC methods in [21], that being that the estimates can be computed in an embarrassingly parallel manner, which provides computational speed-ups versus traditional MCMC algorithms. This is important as our methodology handles two forms of bias, one which arises from the MCMC, and the second arising from the discretization bias associated to a model problem. This is a distinguishment over methods developed by Jacob et al., [21], but also poses improvements, as we will demonstrate later in the paper, over recently developed methods like the SFS method. In particular from our work our highlighting contribution is that we establish that our estimate is unbiased and of finite variance and, either has finite expected cost to compute, or has finite cost with high probability. These results rely heavily upon the works in [19] and [11], where we assume similar assumptions for the former, and the latter of which provides \mathcal{V} -uniform ergodicity of the discretized Markov kernel that we use, where \mathcal{V} is a chosen Lyapunov function. To highlight our theoretical findings we provide numerical experiments on a range of interesting and challenging examples, which arise in statistics and physics. These include a Bayesian logistic regression problem, a double well potential model and finally a Ginzburg-Landau model. We demonstrate both the estimator achieving finite variance and unbiasedness. We compare it to the SFS [38] to demonstrate the gains achieved from our proposed estimator. To demonstrate the significance of our algorithm further, we also compare it to the unbiased Metropolis adjusted Langevin algorithm (U-MALA) on one of our numerical examples.

1.2 Outline

This article is structured as follows. In [Section 2](#) we detail our proposed methodology based on the underdamped Langevin dynamics. This will lead onto [Section 3](#) where we establish that our estimator is unbiased and of finite variance. Numerical experiments are then conducted in [Section 4](#), where we illustrate our methods on several challenging examples. We conclude our remarks and future areas of research in [Section 5](#). Finally the appendix houses a technical result used in [Section 3](#), and the majority of the algorithms introduced.

2 Approach

In this section we first provide a common notation which will be used throughout the manuscript. We then will introduce our approach for unbiased estimation. This will include a review and discussion on the underdamped Langevin dynamics, and in our context. Finally we will discuss

how this is related to our new unbiased estimator in the context of maximal couplings and provide our unbiased estimator through various algorithms.

2.1 Notations

Let $(\mathsf{X}, \mathcal{X})$ be a measurable space. For $\varphi : \mathsf{X} \rightarrow \mathbb{R}$ we write $\mathcal{B}_b(\mathsf{X})$, to denote the collection of bounded measurable functions. For $\varphi \in \mathcal{B}_b(\mathsf{X})$, we write the supremum norm as $\|\varphi\| = \sup_{x \in \mathsf{X}} |\varphi(x)|$. We denote the Borel sets on \mathbb{R}^d as $B(\mathbb{R}^d)$. The d -dimensional Lebesgue measure is written as dx . For a metric $\mathbf{d} : \mathsf{X} \times \mathsf{X} \rightarrow \mathbb{R}^+$ on X and a function $\varphi : \mathsf{X} \rightarrow \mathbb{R}$, $\text{Lip}_{\mathbf{d}}(\mathsf{X})$ are the Lipschitz functions (with finite Lipschitz constants), that is for every $(x, w) \in \mathsf{X} \times \mathsf{X}$, $|\varphi(x) - \varphi(w)| \leq \|\varphi\|_{\text{Lip}} \mathbf{d}(x, w)$. $\mathcal{P}(\mathsf{X})$ denotes the collection of probability measures on $(\mathsf{X}, \mathcal{X})$. For a finite measure μ on $(\mathsf{X}, \mathcal{X})$ and a $\varphi \in \mathcal{B}_b(\mathsf{X})$, the notation $\mu(\varphi) = \int_{\mathsf{X}} \varphi(x) \mu(dx)$ is used. For $(\mathsf{X} \times \mathsf{W}, \mathcal{X} \vee \mathcal{W})$ a measurable space and μ a non-negative finite measure on this space, we use the tensor-product of functions notation for $(\varphi, \psi) \in \mathcal{B}_b(\mathsf{X}) \times \mathcal{B}_b(\mathsf{W})$, $\mu(\varphi \otimes \psi) = \int_{\mathsf{X} \times \mathsf{W}} \varphi(x) \psi(w) \mu(d(x, w))$. Given a Markov kernel $K : \mathsf{X} \rightarrow \mathcal{P}(\mathsf{X})$ and a finite measure μ , we use the notations $\mu K(dx') = \int_{\mathsf{X}} \mu(dx) K(x, dx')$ and $K(\varphi)(x) = \int_{\mathsf{X}} \varphi(x') K(x, dx')$, for $\varphi \in \mathcal{B}_b(\mathsf{X})$. The iterated kernel is $K^n(x_0, dx_n) = \int_{\mathsf{X}^{n-1}} \prod_{i=1}^n K(x_{i-1}, dx_i)$. For $A \in \mathcal{X}$, the indicator function is written as $\mathbb{I}_A(x)$. \mathbb{Z}^+ is the set of non-negative integers. I_d denotes the $d \times d$ identity matrix. The transpose of a vector or matrix x is denoted as x^T . We denote $\min(a, b)$ as $a \wedge b$. The operator ∇ denotes the gradient and Δ the Laplacian, while Δ_l denotes a time step-size of 2^{-l} , $l \in \mathbb{Z}^+$.

2.2 Framework

We consider the case of underdamped Langevin dynamics:

$$dX_t = V_t dt, \tag{2.1}$$

$$dV_t = (b(X_t) - \kappa V_t) dt + \sigma dB_t, \tag{2.2}$$

where $\{B_t\}_{t \geq 0}$ is a standard d -dimensional Brownian motion, $(\kappa, \sigma) \in (0, \infty)^2$ are given friction and diffusion coefficients and $b : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a typically of gradient form, $b = -\nabla U$. In the latter case there is, under fairly weak assumptions an invariant measure of the process $\{X_t, V_t\}_{t \geq 0}$, with Lebesgue density, $\pi(x, v)$:

$$\pi(x, v) \propto \exp \left\{ \frac{-\kappa (2U(x) + \|v\|^2)}{\sigma^2} \right\}.$$

In practice, one often resorts to time-discretization of the dynamics (2.1)-(2.2), for which we focus upon the Euler discretization of step-size $\Delta_l = 2^{-l}$, $l \in \mathbb{N}_0$:

$$X_{(k+1)\Delta_l} = X_{k\Delta_l} + V_{k\Delta_l} \Delta_l, \tag{2.3}$$

$$V_{(k+1)\Delta_l} = V_{k\Delta_l} + (b(X_{k\Delta_l}) - \kappa V_{k\Delta_l}) \Delta_l + \sigma (B_{(k+1)\Delta_l} - B_{k\Delta_l}), \tag{2.4}$$

with $k \in \mathbb{N}_0$. Under assumptions (see [11] Section 2.3), for l large enough, the discrete-time Markov chain expressed as $\{X_{k\Delta_l}, V_{k\Delta_l}\}_{k \in \mathbb{N}_0}$ admits a unique invariant measure η_l and moreover will converge (in an appropriate sense) to it geometrically quickly. In addition, we should have that η_l will converge to π , by the convergence of Euler approximations. As it will facilitate the approach we are to introduce, we shall modify (2.3) to

$$X_{(k+1)\Delta_l} = X_{k\Delta_l} + V_{k\Delta_l} \Delta_l + \sigma_l \Gamma_{k,l}, \tag{2.5}$$

where, $\{\sigma_l\}_{l \in \mathbb{N}_0}$ is any sequence of non-negative and decreasing constants, that converge to zero and for each $l \in \mathbb{N}_0$, $\{\Gamma_{k,l}\}_{k \in \mathbb{N}_0}$ is a sequence of i.i.d. d -dimensional Gaussian random variables of zero mean and covariance matrix the identity multiplied by Δ_l and that this family of sequences are independent of all other random variables (and of each other). As this modification falls under the framework of [11], under assumptions, that for l large enough, the discrete-time Markov chain $\{X_{k\Delta_l}, V_{k\Delta_l}\}_{k \in \mathbb{N}_0}$, following the dynamics (2.4)-(2.5), admits a unique invariant measure π_l (which is likely to be different to η_l) and moreover will converge (in an appropriate sense) to π_l geometrically quickly. In addition, we should have that π_l will converge to π , by the convergence of Euler approximations and that the noise in (2.5) will disappear.

The approach here is as follows: by using only the dynamics (2.4)-(2.5), we will show that one can produce an estimator of $\pi(\varphi \otimes 1) = \int_{\mathbb{R}^{2d}} \varphi(x, v) \pi(x, v) dx dv$ that is unbiased, where $\varphi : \mathbb{R}^{2d} \rightarrow \mathbb{R}$ is π -integrable. In practice, of course, we will only be interested in the marginal on the X -co-ordinate of π . In other words, we will be interested in unbiasedly estimating the integral $\int_{\mathbb{R}^d} \Psi(x) \tilde{\pi}(x) dx$, with $\Psi(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ is $\tilde{\pi}$ -integrable and $\tilde{\pi}(x) \propto \exp\{-(2\kappa/\sigma^2)U(x)\}$ for some differentiable potential function U .

2.3 A Conditional-Type Max-Coupling

In computational statistics a common procedure is to couple various distributions μ_1 and μ_2 . Particular examples of this can include independent coupling, or optimal couplings based on the theory of optimal transport. For this work we make use of maximum couplings [42].

To construct our approach, we will need the following simple idea. Suppose that we have two positive Lebesgue densities $\mu_1(x_1, y)$ and $\mu_2(x_2, y)$ on spaces $\mathbf{X}_1 \times \mathbb{R}^{2d}$ and $\mathbf{X}_2 \times \mathbb{R}^{2d}$, where \mathbf{X}_1 and \mathbf{X}_2 are two possibly different dimensional sub-spaces of powers of the real line. Suppose also, that we know pointwise the conditional densities, for $j \in \{1, 2\}$, $x_j \in \mathbf{X}_j$ fixed:

$$\mu_j(y|x_j) = \frac{\mu_j(x_j, y)}{\int_{\mathbb{R}^d} \mu_j(x_j, y) dy}.$$

Our objective is to sample from a coupling of μ_1 and μ_2 , so that there is a non-zero probability that the Y -co-ordinate can be equal. Let $\tilde{\mu}$ be any coupling of the marginals $\mu_1(x_1)$ and $\mu_2(x_2)$ and $\bar{\mu}(d(y, y')|x_1, x_2)$ be the maximal coupling of $\mu_1(\cdot|x_1)$ and $\mu_2(\cdot|x_2)$ then one way to achieve our given objective is to sample from the distribution associated to the probability:

$$\tilde{\mu}(d(x_1, x_2)) \bar{\mu}(d(y, y')|x_1, x_2),$$

which can be broadly thought of as a posterior. Note that it is straight-forward to show that the marginal of (y, x_j) is μ_j for $j \in \{1, 2\}$. This concept can be straight-forwardly extended to the case of 4 targets, as in [19, Section 3.2.2.], using much the same construction.

2.4 Method

The method that we pursue is an adaptation of the approach in [19] which was originally developed for unbiased inference for Bayesian inverse problems. The contribution here is that we provide an original coupling construction that one might expect is easier to apply than to conventional Markov kernels such as Metropolis-Hastings. In addition, to verify the unbiasedness, one needs to merge

the theory of the underdamped Langevin dynamics [11], and unbiased theory developed in [19]. The former will be discussed later, with various assumptions stated before the main theorem is presented. We will now discuss our unbiased strategy, based on the work in [19].

2.4.1 Overall Strategy

Let $l_* \in \mathbb{N}_0$ be given and \mathbb{P}_L be any positive probability mass function on $\mathbb{N}_{l_*} := \{l_*, l_* + 1, \dots\}$. Let $\{\xi_l\}_{l \in \mathbb{N}_{l_*}}$ be any sequence of independent random variables, such that

$$\begin{aligned}\mathbb{E}[\xi_{l_*}] &= \pi_{l_*}(\varphi) \\ \mathbb{E}[\xi_l] &= \pi_l(\varphi) - \pi_{l-1}(\varphi) =: [\pi_l - \pi_{l-1}](\varphi) \quad l \in \{l_* + 1, l_* + 2, \dots\}.\end{aligned}$$

Now, let L be a random variable with probability \mathbb{P}_L that is independent of the sequence $\{\xi_l\}_{l \in \mathbb{N}_{l_*}}$ then

$$\widehat{\pi}(\varphi) = \frac{\xi_L}{\mathbb{P}_L(L)}, \quad (2.6)$$

is an unbiased estimator of $\pi(\varphi)$; see [29, 35] for the initial statement and proof. Moreover, if

$$\sum_{l \in \mathbb{N}_{l_*}} \frac{\mathbb{E}[\xi_l^2]}{\mathbb{P}_L(l)} < +\infty, \quad (2.7)$$

the estimator $\widehat{\pi}(\varphi)$ has finite variance. There is also the independent sum-estimator, which can work better than this estimator, but we shall not discuss that for now. The main challenge is then to construct the sequence $\{\xi_l\}_{l \in \mathbb{N}_{l_*}}$.

Typically, one will run $M \in \mathbb{N}$ independent replicates of (2.6), where for each replicate i , $l_i \sim \mathbb{P}_L$, and then use the average

$$(\widehat{\pi}(\varphi))_{\text{avg}} := \frac{1}{M} \sum_{i=1}^M (\widehat{\pi}(\varphi))^{(i)},$$

where $(\widehat{\pi}(\varphi))^{(i)}$ represents the i -th independent replicate of the estimate.

2.4.2 Unbiased Approximation of $\pi_{l_*}(\varphi)$

Throughout the section $l \in \mathbb{N}_{l_*}$ is fixed; although we are interested in the case $l = l_*$ the subsequent exposition holds for any fixed l . We will use a strategy that was developed in [16] (see also [21]) for obtaining the given estimator. The approach is to construct a Markov chain on the product space \mathbb{U}^2 , where $\mathbb{U} = \mathbb{R}^{2d}$. From herein, we consider the Markov transition associated to (2.4)-(2.5) over a unit time interval. We denote this kernel as $K_l : \mathbb{U} \rightarrow \mathcal{P}(\mathbb{U})$, where $\mathcal{P}(\mathbb{U})$ are the collection of probability measures on measurable space $(\mathbb{U}, \mathcal{U})$, \mathcal{U} is the Borel σ -field on \mathbb{R}^{2d} . The reason for this, will become apparent as we continue in our exposition.

In order to follow the construction in [16] we will need to construct a Markov chain, $\{U_{n,l}, \tilde{U}_{n,l}\}_{n \in \mathbb{N}_0}$, on $(\mathbb{U}^2, \mathcal{U} \vee \mathcal{U})$ so that marginally, at a given discrete time, new positions of the chain have the

kernel K_l as a marginal. More precisely, we need to construct a Markov kernel, $\check{K}_l : \mathbf{U}^2 \rightarrow \mathcal{P}(\mathbf{U}^2)$, so that for any $(u_l, \tilde{u}_l, A) \in \mathbf{U}^2 \times \mathcal{U}$:

$$\begin{aligned} \int_{A \times \mathbf{U}} \check{K}_l((u_l, \tilde{u}_l), d(u'_l, \tilde{u}'_l)) &= \int_A K_l(u_l, du'_l), \\ \int_{\mathbf{U} \times A} \check{K}_l((u_l, \tilde{u}_l), d(u'_l, \tilde{u}'_l)) &= \int_A K_l(\tilde{u}_l, d\tilde{u}'_l). \end{aligned} \quad (2.8)$$

To build the coupling \check{K}_l we shall decompose this into a mixture of two kernels $\check{Q}_l : \mathbf{U}^2 \rightarrow \mathcal{P}(\mathbf{U}^2)$ and $\check{P}_l : \mathbf{U}^2 \rightarrow \mathcal{P}(\mathbf{U}^2)$. We will explain why this is required, as the discussion progresses. The simulation of the first kernel \check{Q}_l is described in [Algorithm 2](#). On inspection of [Algorithm 2](#), it is clear that this samples from a coupling of K_l which is such that if $u = \tilde{u}$ then $u' = \tilde{u}'$. We note however, that if $u \neq \tilde{u}$, then there is zero probability that $u' = \tilde{u}'$ and being able to achieve the latter is critical in our construction. This motivates our second kernel \check{P}_l whose simulation is given in [Algorithm 3](#).

Remark 2.1. *The kernel \check{P}_l is essentially the same as \check{Q}_l except at the final time step, one has a non-zero probability that $u' = \tilde{u}'$ irrespective of u, \tilde{u} and note that if $u = \tilde{u}$ one will always simulate $u' = \tilde{u}'$. Note also, that it easily follows that \check{Q}_l samples from a coupling of K_l ; the kernel is a simple application of the method in [Section 2.3](#).*

In [Algorithm 3](#), step 4. the maximal coupling can be sampled, for instance, using the maximal coupling algorithm in [\[42\]](#) or the reflection maximal coupling described in [\[21\]](#), where the later performs better for high-dimensional models and therefore we adopt it. Then, for some $\alpha \in (0, 1)$ fixed we set

$$\check{K}_l = \alpha \check{Q}_l + (1 - \alpha) \check{P}_l.$$

The reason for including the kernel \check{Q}_l is that it encourages the positions of $u_{n,l}, \tilde{u}_{n,l}$ to be close, but can never achieve equality; this latter task can be achieved by \check{P}_l and is why this latter kernel is used.

To initialize the Markov chain, for some initial probability $\nu_l \in \mathcal{P}(\mathbf{U})$ we use the following idea as in [\[16\]](#). Let $\tilde{\nu}_l$ be any coupling of (ν_l, ν_l) then the initial probability is taken as:

$$\bar{\nu}_l(d(u', \tilde{u}')) = \int_{\mathbf{U}^2} \tilde{\nu}_l(d(u, \tilde{u})) K_l(u, du') \delta_{\{\tilde{u}\}}(d\tilde{u}'),$$

which evolves from the Markov kernel \check{K}_l and satisfies the properties discussed through [\(2.8\)](#).

Recall we are interested in producing an unbiased estimator of the quantity

$$\pi_l(\varphi) := \varphi(U_{k,l}) + \sum_{n=k+1}^{\infty} \{\varphi(U_{n,l}) - \varphi(U_{n-1,l})\}. \quad (2.9)$$

However doing so can be challenging, in particular because we have an infinite sum. Therefore what can be done instead, as eluded too from the discussion above, is to construct another chain $(\tilde{U}_{n,l})$ such that $U_{n,l} = \tilde{U}_{n,l}$, for $n \geq \tau$, where when the two chains meet, this is known as the *meeting time* τ_l , which is defined as follows

$$\tau_l := \inf\{n \geq 1 : U_{n,l} = \tilde{U}_{n,l}\},$$

where we require that τ_l is almost surely finite and that for each $n \geq \tau_l$, $U_{n,l} = \tilde{U}_{n,l}$ almost surely. Therefore what we can do is replace the infinite limit with the meeting time τ_l with a lag of 1, where now we have the following unbiased estimator of $\pi_l(\varphi)$ for any $k \in \mathbb{N}_0$

$$\widehat{\pi_l(\varphi)}_k := \varphi(U_{k,l}) + \sum_{n=k+1}^{\tau_l-1} \{\varphi(U_{n,l}) - \varphi(\tilde{U}_{n,l})\}, \quad (2.10)$$

with $\sum_{n=k+1}^{\tau_l-1} \{\varphi(U_{n,l}) - \varphi(\tilde{U}_{n,l})\} = 0$ if $\tau_l - 1 < k + 1$.

2.4.3 Assumptions for $\widehat{\pi_l(\varphi)}$

In order to prove that (2.10) is an unbiased estimator, with finite variance and expect cost, we summarize these weak assumptions, without going into exact details.

1. We require *convergence* of the marginal chains, i.e.,

$$\lim_{n \rightarrow \infty} \mathbb{E}[\varphi(U_{n,l})] = \mathbb{E}_\pi[\varphi(U)].$$

2. The meeting time $\tau_l := \inf\{n \geq 1 : U_{n,l} = \tilde{U}_{n,l}\}$ has *geometric tails*

$$\mathbb{P}(\tau > n) \leq C\rho^n, \quad \text{for } C < \infty, \rho \in (0, 1).$$

3. *Faithfulness*: once both chains meet at the meeting time τ_l , then $U_{n,l} = \tilde{U}_{n,l}$ for $n \geq \tau_l$.

We remark that a time-averaged extension is also possible; let $(m, k) \in \mathbb{Z}^+ \times \mathbb{Z}^+$, with $m \geq k$, then we have

$$\widehat{\pi_l(\varphi)}_{T,k,m} := \frac{1}{m-k+1} \sum_{n=k}^m \varphi(U_{n,l}) + \sum_{n=k+1}^{\tau_l-1} \left(1 \wedge \frac{n-k}{m-k+1}\right) \{\varphi(U_{n,l}) - \varphi(\tilde{U}_{n,l})\}, \quad (2.11)$$

which recovers (2.10) in the case $m = k$.

2.4.4 Unbiased Approximation of $[\pi_l - \pi_{l-1}](\varphi)$

Our objective is now to provide, for $l \in \{l_* + 1, l_* + 2, \dots\}$ fixed, an unbiased estimator of $[\pi_l - \pi_{l-1}](\varphi)$, which will essentially use the approach detailed in the previous section. Indeed, one could simply use the method outlined above, independently, for π_l and π_{l-1} and independently for each $l \in \{l_* + 1, l_* + 2, \dots\}$. However, this is unlikely to provide an estimator that can achieve (2.7) and hence the variance of such an approach is infinite and not useful in practice. We therefore present an alternative method.

The idea we use is to generate a Markov chain on (Z, \mathcal{Z}) , where $Z = \mathcal{U}^4$ and $\mathcal{Z} = (\mathcal{U} \vee \mathcal{U}) \vee (\mathcal{U} \vee \mathcal{U})$ and we write

$$Z_{n,l,l-1} = \left((U_{n,l}, \tilde{U}_{n,l}), (U_{n,l-1}, \tilde{U}_{n,l-1}) \right).$$

The associated Markov kernel $\check{K}_{l,l-1} : Z \rightarrow \mathcal{P}(Z)$ will be developed below. The construction of our Markov chain, should be that, marginally $\{U_{n,l}\}_{n \in \mathbb{N}_0}$ (resp. $\{\tilde{U}_{n,l}\}_{n \in \mathbb{N}_0}$) is a Markov chain with

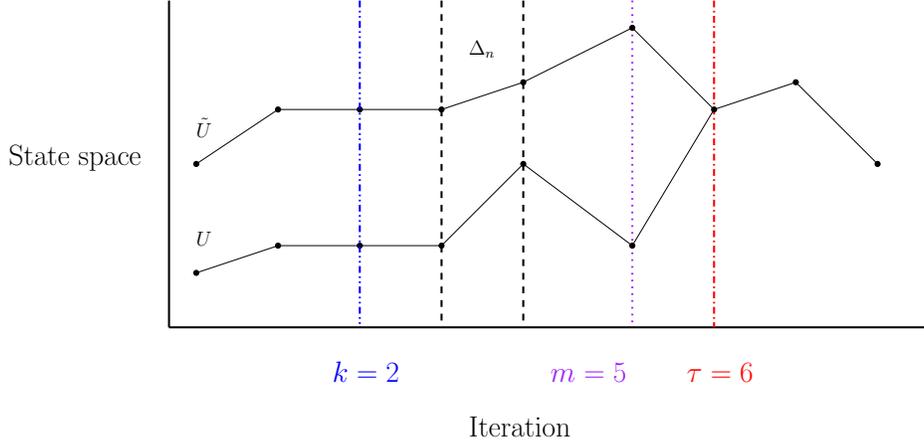


Figure 1: Example illustration of the time-averaged estimator $\widehat{\pi_l(\varphi)}_{T,k,m}$ from (2.11), which includes the meeting time $\tau = 7$ with respect to the two chains U and \tilde{U} . Each line corresponds to a chain, for which we are aiming to couple.

kernel K_l and marginally $\{U_{n,l-1}\}_{n \in \mathbb{N}_0}$ (resp. $\{\tilde{U}_{n,l-1}\}_{n \in \mathbb{N}_0}$) is a Markov chain with kernel K_{l-1} . It is explicitly assumed that (at the very least) $(Z_{n,l,l-1})_{n \in \mathbb{N}_0}$ is constructed so that the stopping time $\tilde{\tau}_{l,l-1} := \tau_l \vee \tau_{l-1}$ is almost surely finite. In addition, the pair of chains on each level should be faithful, i.e. for $s \in \{l, l-1\}$, we have

$$U_{n,s} = \tilde{U}_{n,s}, \text{ for all } n \geq \tau_s. \quad (2.12)$$

Hence for time $n \geq \tilde{\tau}_{l,l-1}$, $Z_{n,l,l-1}$ only has a distinct state on each level.

Our Markov kernel $\check{K}_{l,l-1}$ will be a type of mixture constituting two distinct Markov kernels $\check{Q}_{l,l-1} : Z \rightarrow \mathcal{P}(Z)$ and $\check{P}_{l,l-1} : Z \rightarrow \mathcal{P}(Z)$. We describe the simulation of $\check{Q}_{l,l-1}$ in Algorithm 5. This kernel samples, for $s \in \{l, l-1\}$, from a coupling of K_s which is such that if $u_s = \tilde{u}_s$ then $u'_s = \tilde{u}'_s$ and if $u_s \neq \tilde{u}_s$, then there is zero probability that $u'_s = \tilde{u}'_s$. The simulation is also such that the outputs across the levels should be quite dependent and, in some sense, close. The additional kernel $\check{P}_{l,l-1}$ that we require is described in Algorithm 6. The kernel $\check{P}_{l,l-1}$ is essentially the same as $\check{Q}_{l,l-1}$ except at the final time step, one has a non-zero probability that, for $s \in \{l, l-1\}$, $u'_s = \tilde{u}'_s$ irrespective of u_s, \tilde{u}_s and note that if $u_s = \tilde{u}_s$ one will always simulate $u'_s = \tilde{u}'_s$. Note also, that it easily follows that $\check{Q}_{l,l-1}$ samples from a couplings of K_l and K_{l-1} . In Algorithm 6, step 4. the synchronous pairwise reflection maximal coupling is described in ([19] Section 3.2.3) and the simulation thereof is similar to the reflection maximal coupling. Then for $\alpha \in (0, 1)$ we set

$$\check{K}_{l,l-1} = \mathbb{I}_{D^2}(u_l, \tilde{u}_l, u_{l-1}, \tilde{u}_{l-1}) \check{Q}_{l,l-1} + \mathbb{I}_{(D^2)^c}(u_l, \tilde{u}_l, u_{l-1}, \tilde{u}_{l-1}) [\alpha \check{Q}_{l,l-1} + (1 - \alpha) \check{P}_{l,l-1}],$$

where $D = \{(u, \tilde{u}) \in \mathbb{U}^2 : u = \tilde{u}\}$.

To initialize the Markov chain, for some initial probabilities, $s \in \{l, l-1\}$, $\nu_s \in \mathcal{P}(\mathbb{U})$ we use the following idea. We need a kernel $\bar{K}_{l,l-1} : \mathbb{U}^2 \rightarrow \mathcal{P}(\mathbb{U}^2)$ described in Algorithm 4 for the initialization. Let $\check{\nu}_s$ be any coupling of (ν_s, ν_s) and $\check{\nu}_{l,l-1}$ be a coupling of $(\check{\nu}_l, \check{\nu}_{l-1})$ then the initial probability is taken as:

$$\bar{\nu}_{l,l-1}(dz'_{l,l-1}) = \int_{\mathbb{U}^4} \check{\nu}_{l,l-1}(dz_{l,l-1}) \bar{K}_{l,l-1}((u_l, u_{l-1}), d(u'_l, u'_{l-1})) \delta_{\{\tilde{u}_l, \tilde{u}_{l-1}\}}(d(\tilde{u}'_l, \tilde{u}'_{l-1})).$$

The process is then sampled at subsequent time-points, using the Markov kernel $\check{K}_{l,l-1}$.

Finally then one can estimate $[\pi_l - \pi_{l-1}](\varphi)$ as follows, for any $k \in \mathbb{N}_0$:

$$\widehat{[\pi_l - \pi_{l-1}](\varphi)}_k := \widehat{\pi_l(\varphi)}_k - \widehat{\pi_{l-1}(\varphi)}_k, \quad (2.13)$$

where $\widehat{\pi_s(\varphi)}_k$ is computed using (2.10) One can also employ time-averaging, for $(m, k) \in \mathbb{N}_0 \times \mathbb{N}_0$ satisfying $m \geq k$:

$$\widehat{[\pi_l - \pi_{l-1}](\varphi)}_{T,k,m} := \widehat{\pi_l(\varphi)}_{T,k,m} - \widehat{\pi_{l-1}(\varphi)}_{T,k,m}, \quad (2.14)$$

where $\widehat{\pi_s(\varphi)}_{T,k,m}$ is computed using (2.11)

2.4.5 Final Methodology and Estimator

We now consolidate the above discussion by summarizing our proposed methodology to unbiasedly estimate $\pi(\varphi)$, which is presented in Algorithm 1. We also provide a summary of our proposed methodology, through a simple diagram given in Figure 2.

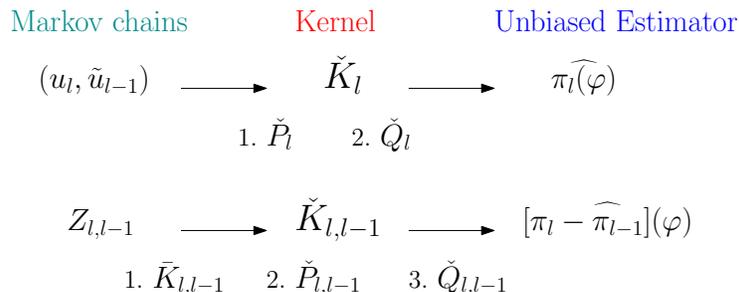


Figure 2: Cartoon description of the required procedure to obtain unbiased estimators, related to the various kernels required.

3 Theoretical Properties

In this section we provide theory related to the discussed methodology in Section 2. Specifically we provide theoretical justification, where we provide our main result which states that our estimator has finite variance, and either has finite expected cost, or has finite cost with a high probability. In order to prove this we require a number of standard assumptions, which have been used in previous works. We omit the proof of theorem in this section, where we defer it to the Appendix. For the below algorithm we have $s \in \{F, C\}$, where F denotes fine, and C denotes coarse, in terms of the level of discretization.

3.1 Assumptions

Before discussing our required assumptions, we provide some new common notation in the setting of Lyapunov functions. Let $0 < C < +\infty$ be a constant and d a metric on \mathbf{U} ; we define the set

$$\mathbf{B}(C, \Delta_l, d) := \{u_{1:4} \in \mathbf{U}^4 : \forall (i, j) \in \{1, \dots, 4\}, d(u_i, u_j) \leq C\Delta_l\}.$$

Algorithm 1 Unbiased estimator $\widehat{\pi(\varphi)}$.

Input: Probability mass function \mathbb{P}_L ,

Initialized Markov chains $(u_0, \tilde{u}_0) = ((x_0, v_0), (\tilde{x}_0, \tilde{v}_0))$, $(u_{0,s}, \tilde{u}_{0,s}) = ((x_{0,s}, v_{0,s}), (\tilde{x}_{0,s}, \tilde{v}_{0,s}))$.

1. Sample $L \sim \mathbb{P}_L$.
2. If $L = l_*$, set $\tilde{U}_{0,l_*} = \tilde{u}_0$ sampling U_{0,l_*} through (2.4) - (2.5) at level l_* up to time 1. Generate $\{U_{n,l}, \tilde{U}_{n,l}\}_{n \in \mathbb{N}}$ according to \check{K}_{l_*} , and compute $\widehat{\pi_0(\varphi)_k}$ in (2.10) or $\widehat{\pi_0(\varphi)}_{T,k,m}$ in (2.11).
3. If $L > l_*$, set $(\tilde{U}_{0,L}, \tilde{U}_{0,L-1}) = (\tilde{u}_{0,F}, \tilde{u}_{0,C})$ sampling $(U_{0,L}, U_{0,L-1})$ from $\bar{K}_{L,L-1}$ by running Algorithm 4. Generate $\{Z_{n,L,L-1}\}_{n \in \mathbb{N}}$ according to $\check{K}_{L,L-1}$, and compute $\widehat{[\pi_L - \pi_{L-1}](\varphi)_k}$ in (2.13) or $\widehat{[\pi_L - \pi_{L-1}](\varphi)}_{T,k,m}$ in (2.14).

Output: Single-term estimator

$$\widehat{\pi(\varphi)}_{S,k} := \frac{1}{\mathbb{P}_L(L)} \left(\mathbb{I}_{\{l_*\}}(L) \widehat{\pi_{l_*}(\varphi)_k} + \mathbb{I}_{\{l_*+1, l_*+2, \dots\}}(L) \widehat{[\pi_L - \pi_{L-1}](\varphi)_k} \right), \quad (2.15)$$

or time-averaged estimator

$$\widehat{\pi(\varphi)}_{S,T,k,m} := \frac{1}{\mathbb{P}_L(L)} \left(\mathbb{I}_{\{l_*\}}(L) \widehat{\pi_{l_*}(\varphi)}_{T,k,m} + \mathbb{I}_{\{l_*+1, l_*+2, \dots\}}(L) \widehat{[\pi_L - \pi_{L-1}](\varphi)}_{T,k,m} \right). \quad (2.16)$$

In what follows, for any $l \in \mathbb{N}_{l_*}$, under assumptions (see [11] Section 2.3) the Markov kernel K_l admits a unique invariant measure π_l . We introduce a Lyapunov function $\mathcal{V} : \mathbf{U} \rightarrow [1, \infty)$ which will be used in our assumptions below. For $f : \mathbf{U} \rightarrow \mathbb{R}$ and a given Lyapunov function \mathcal{V} , we define the collection of functions $\mathcal{L}_{\mathcal{V}} := \{f : \sup_{u \in \mathbf{U}} |f(u)|/\mathcal{V}(u) < +\infty\}$ and we write $\|f\|_{\mathcal{V}} = |f(u)|/\mathcal{V}(u)$.

- (A1)** 1. There exist a $\mathcal{V} : \mathbf{U} \rightarrow [1, \infty)$, $(\lambda, b, \mathbf{C}) \in (0, 1) \times (0, \infty) \times \mathcal{U}$ such that for any $(l, u) \in \mathbb{N}_{l_*} \times \mathbf{U}$

$$K_l(\mathcal{V})(u) \leq \lambda \mathcal{V}(u) + b \mathbb{I}_{\mathbf{C}}(u).$$

2. There exists a $(\rho, C) \in (0, 1) \times (0, \infty)$ such that for any $n \in \mathbb{N}$

$$\sup_{l \in \mathbb{N}_{l_*}} \sup_{u \in \mathbf{U}} \sup_{\varphi \in \mathcal{L}_{\mathcal{V}}} \frac{|K_l^n(\varphi)(u) - \pi_l(\varphi)|}{\mathcal{V}(u)} \leq C \|\varphi\|_{\mathcal{V}} \rho^n,$$

with \mathcal{V} as in 1..

- (A2)** There exist $(C, \rho) \in (0, \infty) \times (0, 1)$ such that:

1. For any $n \in \mathbb{N}$

$$\mathbb{E}[\mathbb{I}_{\{\tau_0 > n\}}] \leq C \rho^n.$$

2. For any $(l, n, z) \in \mathbb{N} \times \mathbb{N} \times \mathsf{X}^4$

$$\mathbb{E}[\mathbb{I}_{\{\tau_l > n\}} | Z_{0,l,l-1} = z] \leq C\rho^n.$$

(A3) There exist a $C < \infty$ and metric on X , $\tilde{\mathbf{d}} : \mathsf{X} \times \mathsf{X} \rightarrow \mathbb{R}^+$ such that for any $(l, \varphi, (x, y)) \in \mathbb{Z}^+ \times \mathcal{B}_b(\mathsf{X}) \cap \text{Lip}_{\tilde{\mathbf{d}}}(\mathsf{X}) \times \mathsf{X} \times \mathsf{X}$

$$|K_l(\varphi)(x) - K_l(\varphi)(y)| \leq C(\|\varphi\| \vee \|\varphi\|_{\text{Lip}})\tilde{\mathbf{d}}(x, y).$$

(A4) There exist $(C, \beta_1) \in (0, \infty) \times (0, \infty)$ such that for any $(l, \varphi) \in \mathbb{N} \times \mathcal{B}_b(\mathsf{X})$

1. $|\pi_l - \pi|(\varphi) \leq C\|\varphi\|\Delta_l^{\beta_1}$.
2. $\sup_{x \in \mathsf{X}} |K_l(\varphi)(x) - K_{l-1}(\varphi)(x)| \leq C\|\varphi\|\Delta_l^{\beta_1}$.

(A5) 1. There exist $(C, \epsilon, \beta_2) \in (0, \infty)^3$, such that for the metric $\tilde{\mathbf{d}}$ as in (A3) and any $(l, n) \in \mathbb{N} \times \mathbb{N}$

$$\begin{aligned} \mathbb{E}[\mathbb{I}_{\mathbb{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c}(Z_{0,l,l-1})] &\leq C\Delta_l^{\beta_2(2+\epsilon)}, \\ \mathbb{E}[\mathbb{I}_{\mathbb{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c \times \mathbb{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})}(Z_{n,l,l-1}, Z_{n-1,l,l-1})] &\leq C\Delta_l^{\beta_2(2+\epsilon)}. \end{aligned}$$

2. $\tilde{d}^{4(2+\epsilon)} \in \mathcal{L}_{\mathcal{V} \otimes \mathcal{V}}$, with $\tilde{\mathbf{d}}, \epsilon$ as in 1. and \mathcal{V} as in (A1).

3.2 Discussion of Assumptions

Let us briefly discuss each assumption in-turn. Assumption (A1) is related to the ergodicity of the underdamped Langevin sampler. The specific result we have used, arises in the work [11], where the authors show that one can attain \mathcal{V} -uniform ergodicity, for the discretized set of equations. These conditions hold so long as l_* is large enough, which we shall also assume. Assumption (A2) has been considered in [17], and is shown to hold in a related context [[17], Lemmata 14 and 20]. Assumption (A3)-(A4) relates to the continuity properties of the kernel and the discretization bias of the problem. In particular (A4) 2. states that moves under pairs of kernels at consecutive levels, stay close on average, given that they are initialized at the same point. Assumption (A5) 1. is a non-standard assumption which can be verified in complex settings, i.e. [[17], Lemma 16]. Also if one can establish that the pairs of chains are uniformly ergodic, with an invariant measure $\tilde{\pi}_{l,l-1}$, then it is reasonable to assume that

$$\tilde{\pi}_{l,l-1} \left((\varphi \otimes 1 - 1 \otimes \varphi)^2 \right) \leq C\Delta_l^{2\beta_2}.$$

The other key assumption which is different in our work is (A5) 2. which is a simple integrability assumption on the metric $\tilde{\mathbf{d}}$ which allows one to remove the need for U to be compact in [19] and is not the case in our context.

3.3 Main Result and Implications

We are now in a position to state our main result, which is that our estimator is unbiased with finite variance.

Proposition 3.1. *Assume (A1-4). Then there exists a choice of positive probability mass function \mathbb{P}_L , such that for the metric \tilde{d} in (A3) and any $\varphi \in \mathcal{B}_b(\mathbf{X}) \cap \text{Lip}_{\tilde{d}}(\mathbf{X})$, (2.15) is an unbiased and finite variance estimator of $\pi(\varphi)$.*

Remark 3.1. *As in [19, Theorem 2.1.], one can easily extend this result to the case of (2.16).*

Proposition 3.1 establishes that there exists a choice of positive probability mass function \mathbb{P}_L which ensures that (2.15) is unbiased and finite variance, but not how to find one. That is the topic of the current discussion and follows very much that given in [19, Section 2.4.3.]. Under our assumptions, noting (2.7), using the argument contained in the Appendix and in [19], the variance of (2.15) is upper-bounded by (with $C < \infty$ a generic constant that does not depend upon l , but whose value may change on each appearance)

$$C \sum_{l \in \mathbb{N}_{l_*}} \frac{\mathbb{E}[\xi_l^2]}{\mathbb{P}_L(l)} \leq C \sum_{l \in \mathbb{N}_{l_*}} \frac{\Delta_l^{2\beta}}{\mathbb{P}_L(l)}, \quad (3.1)$$

with $\xi_l := \widehat{[\pi_l - \pi_{l-1}]}(\varphi)_k$ defined in (2.13) and $\beta = \min\{\beta_1, \frac{\beta_2}{2}\}$. The expected cost is upper-bounded by

$$C \sum_{l \in \mathbb{N}_{l_*}} \mathbb{E}[\check{\tau}_{l,l-1}] \Delta_l^{-1} \mathbb{P}_L(l) \leq C \sum_{l \in \mathbb{N}_{l_*}} \Delta_l^{-1} \mathbb{P}_L(l),$$

where it is assumed that the cost to sample from the kernel K_l is bounded by $C\Delta_l^{-1}$ and that, by A2, one can upper-bound the expected stopping time $\mathbb{E}[\check{\tau}_{l,l-1}]$. If one took $\mathbb{P}_L(l) \propto \Delta_l^\omega$, then if $\omega \in (1, 2\beta)$ both the variance and expected cost are finite. From standard results on Euler discretization of diffusions one might expect that $\beta_1 = 1$ as confirmed in Figure 3, where we estimate the weak error rate for Euler–Maruyama method considering the models in Section 4. We also plot the second moment of the increments ξ_l in Figure 4. If we assume that $\mathbb{E}[\xi_l^2] = \mathcal{O}(\Delta_l^{2\beta})$, then Figure 4 gives an estimate of β for each model, which is approximately 1, and as a result, $\beta_2 > 2$.

4 Numerical Results

In this section we seek to illustrate and verify our theoretical results using three numerical examples arising in statistics and physics. Our experiments will be based on demonstrating both finite variance and unbiasedness of the proposed estimator. We will test the methodology on Bayesian logistic regression problem, a double-well potential and Ginzburg–Landau model. Finally, we will compare our methodology to the unbiased Schrodinger–Föllmer sampler for one particular example.

In all examples below we choose κ and σ such that $2\kappa = \sigma^2$ and set b in (2.4) to $b(x) = -\nabla U(x)$, where $U = -\log \pi$ and π is the density of interest that we aim to sample from. The objective is to estimate, unbiasedly, the expectation $\pi(\varphi)$. We find that the algorithm works quite well when σ is assigned a large value, but it may collapse for small values. In our simulations, we set $\sigma = 3$ and take $\varphi(x) = x$.

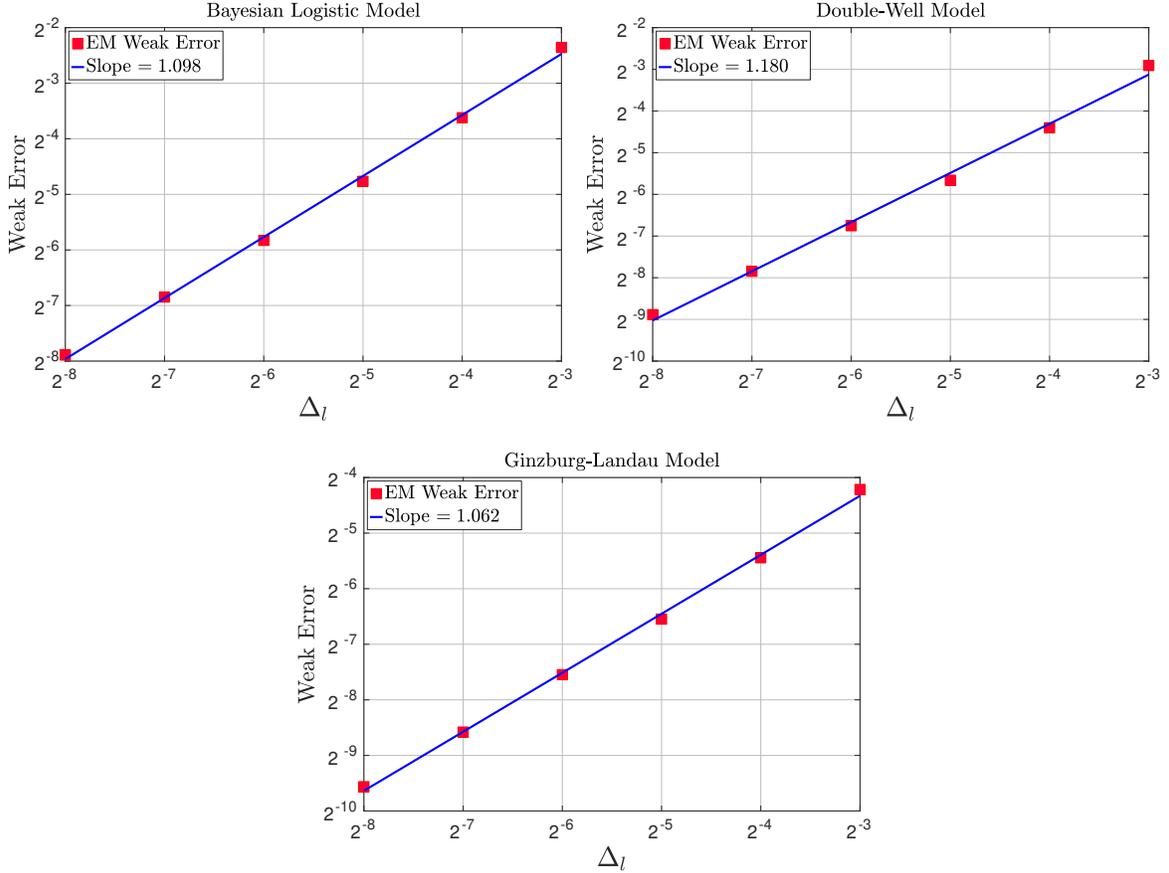


Figure 3: Euler–Maruyama (E–M) weak error $|\pi_l - \pi|(\varphi)$ against the time step-size Δ_l for the models in Section 4, with l the discretization level, $\varphi(u) = u_1$ (the first co-ordinate), $\pi(\varphi)$ the reference expectation which we take to be the average of 5200 samples of Euler-discretized $\pi_L(\varphi)$ at level $L = 18$, a very high level of discretization. These plots verify that the weak error rate β_1 of E–M is approximately 1 as expected.

4.1 Bayesian Logistic Regression

We consider the binary logistic regression model in which the binary observations $\{Y_i\}_{i=1}^n$ are conditionally independent Bernoulli random variables such that $Y_i \in \{0, 1\}$ and

$$\mathbb{P}(Y_i = 1 | X_i = x_i, \alpha) = \rho(\alpha^T x_i),$$

where $\rho : \mathbb{R} \rightarrow (0, 1)$ defined by $\rho(w) = e^w / (1 + e^w)$ is the logistic function and X_i , $\alpha \in \mathbb{R}^d$ are the covariates and the unknown regression coefficients, respectively. The prior density for the parameter α is a multivariate normal Gaussian given by

$$pr(\alpha) = \phi(\alpha; 0, \Sigma_\alpha),$$

where Σ_α is defined through its inverse $\Sigma_\alpha^{-1} = \frac{1}{n}(\sum_{i=1}^n X_i X_i^T)$. The covariate vectors $\{X_i\}_{i=1}^n$ are sampled independently from $\mathcal{U}\{-1, 1\}^d$ which then are standardized. The density of the posterior distribution of α is then given by

$$\pi(\alpha | \{X_i = x_i, Y_i = y_i\}_{i=1}^n) \propto \exp(-U(\alpha)),$$

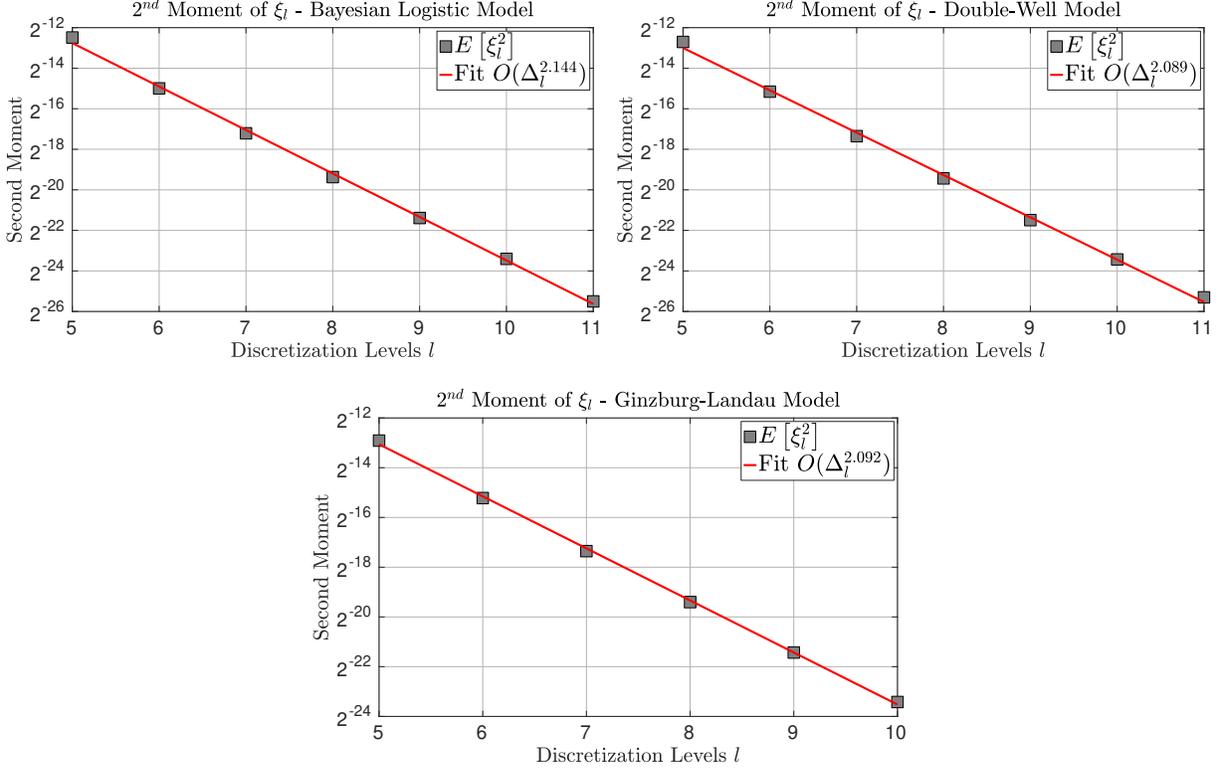


Figure 4: Second moment of the time-averaged estimator increments $\xi_l = \widehat{[\pi_l - \pi_{l-1}](\varphi)}_{S,T,k,m}$ against discretization level l for the different models in Section 4. The function considered here is $\varphi(u) = u_1$, $u \in \mathbb{R}^{2d}$. These plots provide a numerical estimate of the parameter β in (3.1) of value 1.

where

$$U(\alpha) = - \sum_{i=1}^n [y_i \alpha^T x_i - \log(1 + \exp(\alpha^T x_i))] + \frac{1}{2} \alpha^T \Sigma_\alpha^{-1} \alpha,$$

and

$$\nabla U(\alpha) = - \sum_{i=1}^n y_i x_i + \sum_{i=1}^n \frac{\exp(\alpha^T x_i) x_i}{1 + \exp(\alpha^T x_i)} + \alpha^T \Sigma_\alpha^{-1}.$$

We take $d = 5$ and $n = 100$. The reference expectation is the mean of 2.6×10^8 samples computed through the unbiased MCMC method proposed in [21] with a random-walk Metropolis-Hasting kernel.

4.2 Double-Well Model

For our second example we consider sampling from $\pi(x) = \exp(-U(x))$, where U is a double-well potential given by

$$U(x) = \frac{1}{4} \|x\|_2^4 - \frac{1}{2} \|x\|_2^2, \quad x \in \mathbb{R}^d. \quad (4.1)$$

The double-well potential is one of several quartic potentials of substantial importance in quantum mechanics and quantum field theory [25] for the investigation of different physical phenomena or

mathematical features. The gradient of the potential is given by

$$\nabla U(x) = (\|x\|_2^2 - 1) x.$$

In our simulations, we take $d = 100$. The true reference expectation is $\pi(\varphi) = 0$.

4.3 Ginzburg-Landau Model

In the final example we consider the Ginzburg-Landau (GL) model [15, 20, 27], a model that describes a thermodynamic system that undergoes continuous phase transitions at a temperature $T = T_c$, from a high-temperature symmetric phase to a low-temperature ordered phase in which some symmetry is broken. We denote by the random variable $\psi \in \mathbb{R}$ the order parameter, which is assumed to be spatially dependent, i.e., $\psi = \psi(x)$, $x \in \Omega \subseteq \mathbb{R}^3$, and $\nabla_x \psi \cdot n = 0$, where n is a unit vector normal to the boundary $\partial\Omega$. In the absence of external fields, the probability of a fluctuation $\psi(x)$ is given by

$$\pi(\psi) \propto \exp \{-U(\psi(x))\},$$

with U is the GL free energy functional defined via

$$U(\psi(x)) = \int_{\Omega} \left[\frac{1 - \bar{T}}{2} \psi^2(x) + \frac{\gamma \bar{T}}{2} \|\nabla_x \psi(x)\|^2 + \frac{\zeta \bar{T}}{4} \psi^4(x) \right] dx,$$

$\bar{T} = T_c/T$ and $\gamma, \zeta > 0$. We now consider the discretized GL model on a 3d-lattice with $d = d_0^3$ sites, $(x_{i,j,k})_{i,j,k=1}^{d_0}$, where $d_0 \in \mathbb{N}$, with spacing equal to one. Periodic boundary conditions are applied to the system where we have $x_{11,j,k} = x_{1,j,k}$ and $x_{0,j,k} = x_{10,j,k}$ and similar situation for the second and the third coordinates. Each site represents a random variable $\psi_{i,j,k} \in \mathbb{R}$. Set $\boldsymbol{\psi} := (\psi_{i,j,k})_{i,j,k=1}^{d_0} \in \mathbb{R}^d$, then an approximation to π is given by

$$\tilde{\pi}(\boldsymbol{\psi}) \propto \exp \{-\tilde{U}(\boldsymbol{\psi})\},$$

with

$$\tilde{U}(\boldsymbol{\psi}) = \sum_{i,j,k=1}^{d_0} \left[\frac{1 - \bar{T}}{2} \psi_{i,j,k}^2 + \frac{\gamma \bar{T}}{2} ((\psi_{i+1,j,k} - \psi_{i,j,k})^2 + (\psi_{i,j+1,k} - \psi_{i,j,k})^2 + (\psi_{i,j,k+1} - \psi_{i,j,k})^2) + \frac{\zeta \bar{T}}{4} \psi_{i,j,k}^4 \right],$$

such that a forward finite difference is used to approximate $\nabla_x \psi$. The functional derivative $\delta U(\psi(x'))/\delta \psi(x)$

is given by

$$\begin{aligned} \frac{\delta U(\psi(x'))}{\delta \psi(x)} &= \int_{\Omega} \left[(1 - \bar{T}) \psi(x') \delta(x - x') + \gamma \bar{T} \nabla_{x'} \psi(x') \cdot \nabla_{x'} \delta(x - x') + \zeta \bar{T} \psi^3(x') \delta(x - x') \right] dx' \\ &= (1 - \bar{T}) \psi(x) - \gamma \bar{T} \Delta_x \psi(x) + \zeta \bar{T} \psi^3(x), \end{aligned} \quad (4.2)$$

where we used the following identities

$$\begin{aligned} \frac{\delta}{\delta \psi(x)} \psi(x') &= \delta(x - x'), \\ \int_{\Omega} \nabla_{x'} \psi(x') \cdot \nabla_{x'} \delta(x - x') dx' &= - \int_{\Omega} \Delta_{x'} \psi(x') \delta(x - x') dx' = -\Delta_x \psi(x). \end{aligned}$$

Here $\delta(x)$ denotes the Dirac delta function. In the last identity we used integration by parts with the fact that $\nabla_x \psi \cdot n = 0$ for any normal vector n on the boundary. On the lattice, an approximation to the functional derivative in (4.2), which we denote by $\widetilde{\nabla_\psi U}$, is given by

$$\widetilde{\nabla_\psi U} = \left((1 - \bar{T})\psi_{i,j,k} + \gamma \bar{T} \vartheta_{i,j,k} + \zeta \bar{T} \psi_{i,j,k}^3 \right)_{i,j,k=1}^{d_0},$$

and

$$\vartheta_{i,j,k} = \{6\psi_{i,j,k} - (\psi_{i+1,j,k} + \psi_{i-1,j,k} + \psi_{i,j+1,k} + \psi_{i,j-1,k} + \psi_{i,j,k+1} + \psi_{i,j,k-1})\},$$

where a second-order central finite difference is used to approximate $\Delta_x \psi$. Then $b - \widetilde{\nabla_\psi U}$ in (2.4). The true reference expectation is $\pi(\varphi) = 0 \in \mathbb{R}^d$, where $d = 10^3$. We also set $\bar{T} = 2$, $\gamma = 0.1$ and $\zeta = 0.5$.

4.4 Simulation Settings

The time-averaged, unbiased estimator

$$(\widehat{\pi}(\varphi)_{S,T,k,m})_{\text{avg}} := \frac{1}{M} \sum_{i=1}^M (\widehat{\pi}(\varphi)_{S,T,k,m})^{(i)}, \quad (4.3)$$

where $(\widehat{\pi}(\varphi)_{S,T,k,m})^{(i)}$ is the estimator presented in (2.16) in Algorithm 1 and the cost of the single-level, time-averaged estimator

$$(\widehat{\pi}_L(\varphi)_{T,k,m})_{\text{avg}} := \frac{1}{M} \sum_{i=1}^M (\widehat{\pi}_L(\varphi)_{T,k,m})^{(i)}, \quad (4.4)$$

where $(\widehat{\pi}_L(\varphi)_{T,k,m})^{(i)}$ is presented in (2.11) and L is the discretization level. We will compare the cost of both estimators versus the mean-squared errors (MSEs) that are obtained by running 50 independent simulations of each estimator and are given by

$$\text{MSE}_{ub} = \frac{1}{50} \sum_{j=1}^{50} \left[(\widehat{\pi}(\varphi)_{S,T,k,m})_{\text{avg}}^{(j)} - \pi(\varphi) \right]^2, \quad \text{and} \quad \text{MSE}_s = \frac{1}{50} \sum_{i=1}^{50} \left[(\widehat{\pi}_L(\varphi)_{T,k,m})_{\text{avg}}^{(j)} - \pi(\varphi) \right]^2,$$

where $\pi(\varphi)$ is the reference expectation. Let τ be the stopping time, we set $k = 100$ and $m = \min\{2k, \tau - 1\}$, however, if $\tau < k - 1$, we set $k = 0.5 \tau$.

For a given $\epsilon > 0$, the goal is to obtain an MSE of order ϵ^2 . Therefore, for the single-level estimator, we set $L = \mathcal{O}(-\log_2(\epsilon))$ and $M = \mathcal{O}(\epsilon^{-2})$. For the unbiased estimator, in practice, one has to truncate the values of L . As a result, we set $\mathbb{P}_L(l) \propto 2^{-1.5l} \mathbb{I}_{\{l_*, \dots, l_{\max}\}}(l)$, where $l_*, l_{\max} \in \mathbb{N}_0$ and $l_* < l_{\max}$. With this choice, the probability of simulating (2.4)-(2.5) at a high discretization level is very small. In all examples, we take $l_* = 5$, $l_{\max} = 12$ and $M = \mathcal{O}(\epsilon^{-2})$ as in the single-level estimator. The cost to compute the single-level estimator is $\mathcal{C}_s := 2^{L+1} \sum_{j=1}^M \tau_L^{(j)} = 2^{L+1} M \max_j \{\tau_L^{(j)}\} \leq \mathbf{C} 2^L M = \mathcal{O}(\epsilon^{-3})$ where $\tau_L^{(j)}$ is the stopping time for the j -th replicate that is assumed to be bounded by A2. The cost of the unbiased estimator is $\mathcal{C}_{ub} := \sum_{j=1}^M \mathcal{C}^{(j)}$ where

$$\mathcal{C}^{(j)} = \begin{cases} \tau_{L_j} 2^{L_j+1} & \text{if } L_j = l_*, \\ \check{\tau}_{L_j, L_{j-1}} (2^{L_j+1} + 2^{L_j}) & \text{if } L_j > l_*, \end{cases} \quad L_j \sim \mathbb{P}_L.$$

Remark 4.1. We note that even though the rate for the single-level estimator should be $\mathcal{O}(\epsilon^{-3})$, what we observe numerically is different. As we will discuss in the next subsection, the rates obtained are worse, where it seems that we attain a rate of $\mathcal{O}(\epsilon^{-(3+\delta)})$, for $\delta > 0$. This is for each individual numerical example and we expect that this is associated to the drift coefficient in the dynamics, which are in-turn affected by the target. The regularity of the drift will influence the practical rate of convergence of the discretization scheme and the time convergence of the dynamics.

By [Remark 3.1](#), the expected cost of each replicate is bounded by a constant C and hence the overall expected cost is $\overline{C}_{\text{ub}} \leq CM = \mathcal{O}(\epsilon^{-2})$. Notice that in fact the cost of sampling from \check{Q}_L is 2^L and the cost of sampling from \check{P}_L is $2^L - 1$ plus the cost of sampling from a maximal coupling, which we take it to be one, and therefore the cost of sampling from $\check{K}_L = \alpha\check{Q}_L + (1 - \alpha)\check{P}_L$ is 2^{L+1} .

We will also compare the unbiased SFS algorithm [\[38\]](#) with our unbiased ULD sampler [Algorithm 1](#). For a fair comparison, we consider comparing the average machine time (in seconds) that one independent realization of the 50 realizations takes on a workstation of 52 cores. The unbiased SFS estimator is based on the SF diffusion process,

$$dX_t = b(X_t, t)dt + dW_t, \quad X_0 = 0, \quad t \in [0, 1],$$

where the drift term has the specific form

$$b(x, t) = \nabla \log \mathbb{E}[f(x + W_{1-t})] = \frac{\mathbb{E}_\phi[\nabla f(x + \sqrt{1-t}Z)]}{\mathbb{E}_\phi[f(x + \sqrt{1-t}Z)]} = \frac{1}{\sqrt{1-t}} \frac{\mathbb{E}_\phi[Zf(x + \sqrt{1-t}Z)]}{\mathbb{E}_\phi[f(x + \sqrt{1-t}Z)]},$$

and f corresponds to the analogous of a likelihood function for a standard Gaussian prior, i.e. for $z \in \mathbb{R}^d$, $f(z) = \pi(z)/\phi(z)$, with $\phi(z)$ the standard d -dimensional Gaussian density. For $l \in \{0, 1, \dots\}$ and $k \in \{0, 1, \dots, \Delta_l^{-1} - 1\}$, we discrete the above SDE as

$$\widetilde{X}_{(k+1)\Delta_l}^{l,N} = \widetilde{X}_{k\Delta_l}^{l,N} + \hat{b}(\widetilde{X}_{k\Delta_l}^{l,N}, k\Delta_l)\Delta_l + W_{(k+1)\Delta_l} - W_{k\Delta_l},$$

where, for $x \in \mathbb{R}^d$

$$\hat{b}(x, k\Delta_l) = \frac{\frac{1}{N} \sum_{i=1}^N \nabla f(x + \sqrt{1-k\Delta_l}Z^i)}{\frac{1}{N} \sum_{i=1}^N f(x + \sqrt{1-k\Delta_l}Z^i)},$$

and, for $i \in \{1, \dots, N\}$, $Z^i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_d(0, I)$. The idea of the unbiased SFS is based upon double-randomization techniques (e.g. [\[22, 23\]](#)) where both N and l are chosen at random in a specific manner.

4.5 Results and Discussion

In [Figure 5](#) (a)–(c), we plot the cost of the estimators in [\(4.3\)](#)–[\(4.4\)](#) against their MSE for the models above. The MSE of the proposed unbiased estimator presented in [Algorithm 1](#) decays at the optimal rate of $1/\overline{C}_{\text{ub}} = \mathcal{O}(\epsilon^2)$ as shown in the plots, which is as expected, and outperforms that of the single-level estimator given in [\(4.4\)](#). As eluded to, from [Remark 4.1](#) the actual rate which we obtain for each model is around $\mathcal{O}(\epsilon^{-(3+\delta)})$, for $\delta > 0$, which differs for each model. This is especially the case for the GL model which, compared to the other models, has a more complicated drift term. Therefore we believe this constitutes to a modified rate for the single-level

estimator. Further investigation is needed to verify this, but this is beyond the scope of this work and we leave it for future work. Nonetheless, as stated, the unbiased estimator outperforms the single estimator for each model example.

In Figure 5 (d), the cost of our unbiased estimator and that of the SFS, measured in seconds, is plotted against the MSE of each method. As we can see, the new methodology beats the unbiased SFS and in fact the MSE decays faster than expected by the theory. This suggesting to attain similar values of MSE, our unbiased estimator takes considerably less time than that using the SFS, as presented in [38].

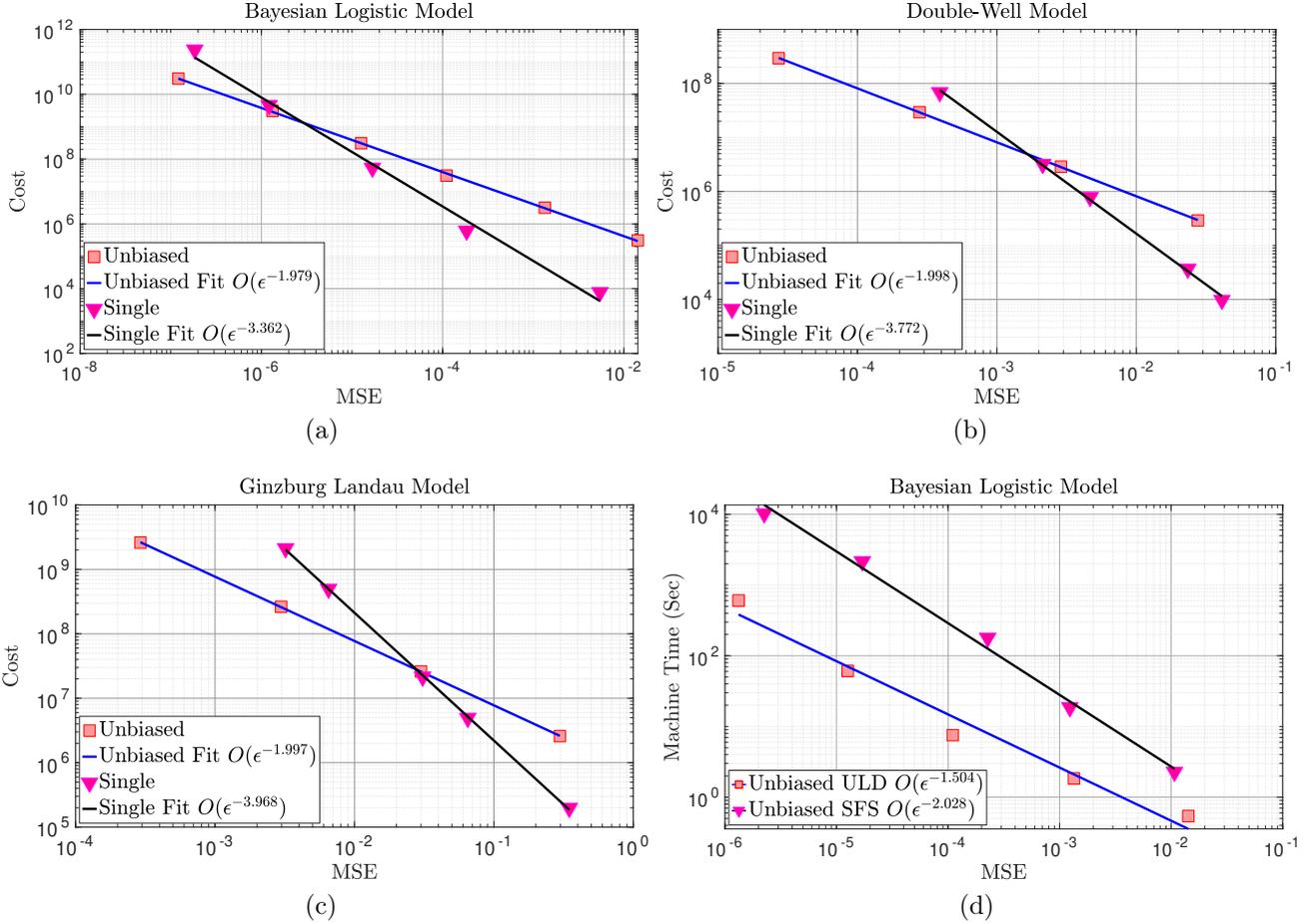


Figure 5: In (a)–(c), we plot the cost versus MSE for the single-level, time-averaged estimator in (2.11) and the time-averaged unbiased estimator presented in (2.16) in Algorithm 1. In (d), we plot the machine time versus MSE for the unbiased estimator in Algorithm 1 and the unbiased SFS estimator in [38] for the Bayesian logistic model.

4.6 Comparison to Unbiased MALA

Our final numerical experiment we present in this work is a comparison of our method to unbiased methodologies proposed by Jacob et al. [21]. In particular we will compare it to the unbiased

MALA, which is a well-known Metropolis Hastings (MH) method based on the following Langevin dynamics

$$dX_t = -\frac{1}{2}\nabla \log \pi(X_t)dt + \sqrt{2\delta^{-1}}dB_t,$$

with discretization

$$X_{(k+1)\Delta_t} = X_{k\Delta_t} - \frac{\Delta_t}{2}\nabla \log \pi(X_{k\Delta_t}) + \sqrt{2\Delta_t\delta^{-1}}(B_{(k+1)\Delta_t} - B_{k\Delta_t}),$$

where where $\{B_t\}_{t \geq 0}$ is a standard d -dimensional Brownian motion and $\delta > 0$ denotes the inverse temperature. The acceptance probability associated with it arises from the usual MH-type algorithms. Further details on MALA can be found in the following references [32, 36, 37]. Our numerical example will consist of a comparison between our unbiased scheme and U-MALA, which is tested on the double-well model (4.1). The U-MALA was first discussed in the work of Heng et al. [18] as a simplified version of the HMC coupling. As a result the coupling associated with U-MALA is much simpler than the U-HMC couplings, which follow from [21], which also exploit synchronous maximal couplings. Our choice of using this model, is that the density of interest is bimodal, which should constitute to a difference in performance over the toy logistic regression model. The dimension is chosen as $d = 75$ and again we run 50 independent simulations to compute the MSE. The other parameter choices are consistent with that discussed in Section 4.4. We also set $\delta = 1$ for our experiments.

We present our simulations in Figure 6. As we can observe for both subplots the MSE-to-cost ratio is better, related to both CPU and theoretical cost, for Algorithm 1. The ratio of the rates between each unbiased estimators is consistent in both subplots. As mentioned we believe the reasons for the under-performance of the unbiased MALA proposed by [18, 21] are due to bimodality of $U(x)$ and the fact the unbiasedness is only related to the bias coming from the MCMC, not the discretization error, whereas our algorithm aims at removing both the MCMC and discretization bias.

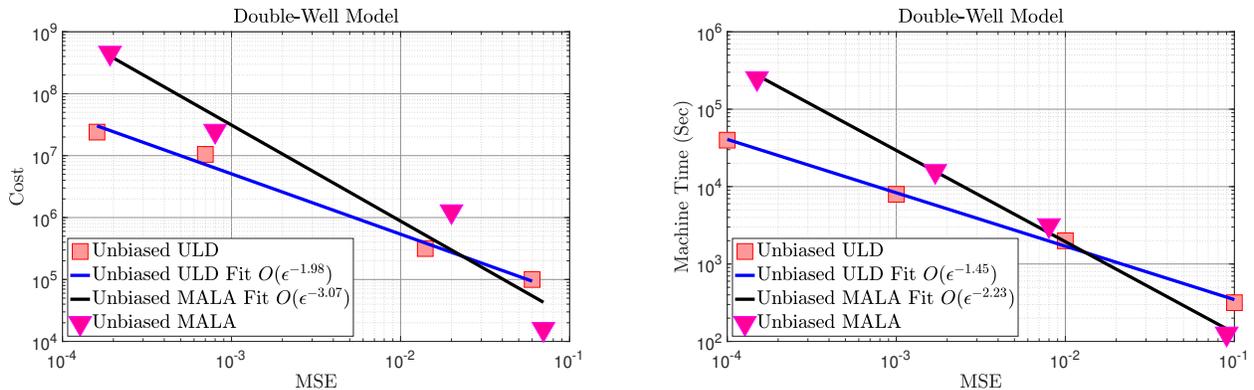


Figure 6: This plot compares the time-average unbiased estimator (2.16) in Algorithm 1 to the unbiased MALA. Left: Theoretical cost vs MSE. Right: CPU time cost vs MSE.

5 Conclusion

Our motivation from this work was to provide a new unbiased estimator, which is based on the discretized underdamped Langevin dynamics (ULD) (2.3)-(2.4). The ULD has sparked recent interest in both the statistics and machine learning community, and as a result we wanted to see if such a dynamics could be used in the context of unbiased estimation. We introduced a new methodology based on ULD, which was based on the double-randomization schemes used for unbiased estimation. Subsequently we proved that our new estimator is unbiased with finite variance under suitable assumptions. To verify our theory, we implemented our methodology on a range of interesting model problems, such as the stochastic Ginzburg-Landau model, the double-well model and a Bayesian logistic regression problem. We also justified such an estimator by comparing it to our known estimators, developed in a similar manner, such as the unbiased Schrödinger–Föllmer sampler presented in [38] and the unbiased Metropolis adjusted Langevin algorithm presented in [18, 21].

From this work, there are a number of research directions one can consider. A rather obvious one would be to use the current methodology aimed at unbiased estimation for the both the score function and the Hessian [6, 17]. This has already been considered previously, but in the context of our methodology which may prove to be more useful. Another direction is to extend the ergodicity result presented in [11], where one does not have the requirement that l_* has to be large enough to ensure \mathcal{V} -uniform ergodicity. This may prove to be challenging, but from a theoretical perspective would be of interest, especially if one can attain geometric ergodicity. One could also consider the mean-field ULD, which has shown to be promising for deep neural architectures [26], in terms of the trainability of two-layer neural networks [8, 28]. Related to the above point, another extension could be the perturbed ULD [10], which has demonstrated improvements for sampling even more complex probability measures. A final direction could be providing an extension related to the work of [45]. The authors were able to provide an upper bound on one-step meeting probabilities, related to the proposal and the acceptance probability based on particular setup.

Acknowledgements

All three authors were supported by KAUST baseline funding. We would like to thank the editors and reviewers for their guidance which has greatly improved the article.

A Proof

The proof of Proposition [Proposition 3.1](#) is virtually identical to that of [19, Theorem 2.1]. There is only place where the proof has to be modified and we give the result below. Below C is a generic constant which does not depend upon n, l and whose value changes upon each appearance. The expectation operator \mathbb{E} relates to law which generates the simulated process used to compute (2.15).

Lemma A.1. *Assume (A1,4). Then there exists a $C \in (0, \infty)$ such that for any $(l, n) \in \mathbb{N}_{l_*} \times \mathbb{Z}^+$*

we have:

$$\begin{aligned} \mathbb{E}[\mathbb{I}_{\mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c}(Z_{n,l,l-1})] &\leq C(n+1)\Delta_l^{\beta_2(2+\epsilon)}, \\ \max \left\{ \mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon}], \mathbb{E}[\tilde{\mathbf{d}}(\tilde{U}_{n,l}, \tilde{U}_{n,l-1})^{2+\epsilon}] \right\} &\leq C(n+1)\Delta_l^{\beta_2(1+\frac{\epsilon}{2})}, \end{aligned}$$

where β_2 and ϵ are as (A4).

Proof. The first inequality is proved in an identical manner to [19, Lemma A.3.], so we consider only the second inequality. This latter inequality is shown for $\mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon}]$, only, as the argument for the other term is the same up to changes in notation. Let $n \in \mathbb{N}$, then we have

$$\mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon}] = \mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon} \mathbb{I}_{\mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c}(Z_{n,l,l-1})] + \mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon} \mathbb{I}_{\mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})}(Z_{n,l,l-1})].$$

For the second term on the R.H.S. as $Z_{n,l,l-1} \in \mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})$, we have $\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon} \leq C\Delta_l^{\beta_2(2+\epsilon)}$, so we focus on the first term on the R.H.S.. Applying Cauchy-Schwarz and the first-statement of the Lemma, we have

$$\mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon} \mathbb{I}_{\mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c}(Z_{n,l,l-1})] \leq \mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2(2+\epsilon)}]^{1/2} C(n+1)\Delta_l^{\beta_2(1+\frac{\epsilon}{2})}.$$

Then using $\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{4(2+\epsilon)} \in \mathcal{L}_{\mathcal{V} \otimes \mathcal{V}}$, we get

$$\mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon} \mathbb{I}_{\mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c}(Z_{n,l,l-1})] \leq \mathbb{E}[\{\mathcal{V}(U_{n,l})\mathcal{V}(U_{n,l-1})\}^{\frac{1}{2}}]^{1/2} C(n+1)\Delta_l^{\beta_2(1+\frac{\epsilon}{2})}.$$

Then applying Cauchy-Schwarz and using (A1) 1. we have

$$\mathbb{E}[\tilde{\mathbf{d}}(U_{n,l}, U_{n,l-1})^{2+\epsilon} \mathbb{I}_{\mathbf{B}(C, \Delta_l^{\beta_2}, \tilde{\mathbf{d}})^c}(Z_{n,l,l-1})] \leq C(n+1)\Delta_l^{\beta_2(1+\frac{\epsilon}{2})},$$

from which the proof can be completed. \square

B Algorithms

In this Appendix we provide each algorithm required for our unbiased estimator. This is related to the simulation from the associated kernels $\check{Q}_l, \check{Q}_{l,l-1}, \check{P}_l, \check{P}_{l,l-1}$, which are discussed in Algorithms 2 to 6. For the earlier algorithms, fuller details can be found in [42]. Recall that our objective is to construct the coupling, or coupled kernels, \check{K}_l and $\check{K}_{l,l-1}$, which can be decomposed with the abovely stated kernels where

$$\begin{aligned} \check{K}_l &= \alpha\check{Q}_l + (1-\alpha)\check{P}_l, \\ \check{K}_{l,l-1} &= \mathbb{I}_{D^2}(u_l, \tilde{u}_l, u_{l-1}, \tilde{u}_{l-1})\check{Q}_{l,l-1} + \mathbb{I}_{(D^2)^c}(u_l, \tilde{u}_l, u_{l-1}, \tilde{u}_{l-1})[\alpha\check{Q}_{l,l-1} + (1-\alpha)\check{P}_{l,l-1}], \end{aligned}$$

where $\alpha \in (0, 1)$. Finally we require the coupling $\overline{K}_{l,l-1}$, which is required for the initialization before sampling from $\check{K}_{l,l-1}$, and is presented in Algorithm 4.

Algorithm 2 : Sampling from kernel \check{Q}_l .

1. **Input:** $l, (u, \tilde{u}) = ((x_0, v_0), (\tilde{x}_0, \tilde{v}_0))$.

2. Sample $(\Gamma_{0,l}, \dots, \Gamma_{\Delta_l^{-1}-1,l})$ and $(B_{\Delta_l}, B_{2\Delta_l} - B_{\Delta_l}, \dots, B_1 - B_{1-\Delta_l})$.
 3. Run the recursion (2.4)-(2.5) with $\{\Gamma_{k,l}\}_{k=0}^{\Delta_l^{-1}-1}$ and $\{B_{(k+1)\Delta_l} - B_{k\Delta_l}\}_{k=0}^{\Delta_l^{-1}-1}$ up-to time 1.
 4. **Output:** $u' = (x_1, v_1)$ and $\tilde{u}' = (\tilde{x}_1, \tilde{v}_1)$ as generated in step 3.
-

Algorithm 3 : Sampling from kernel \tilde{P}_l .

1. **Input:** $l, (u, \tilde{u}) = ((x_0, v_0), (\tilde{x}_0, \tilde{v}_0))$.
 2. Sample $(\Gamma_{0,l}, \dots, \Gamma_{\Delta_l^{-1}-2,l})$ and $(B_{\Delta_l}, B_{2\Delta_l} - B_{\Delta_l}, \dots, B_{1-\Delta_l} - B_{1-2\Delta_l})$.
 3. Run the recursion (2.4)-(2.5) with $\{\Gamma_{k,l}\}_{k=0}^{\Delta_l^{-1}-2}$ and $\{B_{(k+1)\Delta_l} - B_{k\Delta_l}\}_{k=0}^{\Delta_l^{-1}-2}$ up-to time $1 - \Delta_l$.
 4. Sample $((X_1, V_1), (\tilde{X}_1, \tilde{V}_1)) | ((x_{1-\Delta_l}, v_{1-\Delta_l}), (\tilde{x}_{1-\Delta_l}, \tilde{v}_{1-\Delta_l}))$ from a maximal coupling of $p_l(x_1, v_1 | x_{1-\Delta_l}, v_{1-\Delta_l})$ and $p_l(\tilde{x}_1, \tilde{v}_1 | \tilde{x}_{1-\Delta_l}, \tilde{v}_{1-\Delta_l})$, where $p_l \sim N_2(m, C)$ where m, C are determined from (2.4)-(2.5).
 5. **Output:** $u' = (x_1, v_1)$ and $\tilde{u}' = (\tilde{x}_1, \tilde{v}_1)$ as generated in step 4..
-

Algorithm 4 : Sampling from coupled kernel $\overline{K}_{l,l-1}$.

1. **Input:** $l \in \{l_* + 1, l_* + 2, \dots\}$, and $(u_l, u_{l-1}) = ((x_{0,l}, v_{0,l}), (x_{0,l-1}, v_{0,l-1}))$.
 2. Sample $(\Gamma_{0,l}, \dots, \Gamma_{\Delta_l^{-1}-1,l})$ and $(B_{\Delta_l}, B_{2\Delta_l} - B_{\Delta_l}, \dots, B_1 - B_{1-\Delta_l})$.
Concatenate to obtain $(\Gamma_{0,l-1}, \dots, \Gamma_{\Delta_{l-1}^{-1}-1,l-1})$
and $(B_{\Delta_{l-1}}, B_{2\Delta_{l-1}} - B_{\Delta_{l-1}}, \dots, B_1 - B_{1-\Delta_{l-1}})$.
 3. For $s \in \{l, l-1\}$: run the recursion (2.4)-(2.5) with $\{\Gamma_{k,s}\}_{k=0}^{\Delta_s^{-1}-1}$ and $\{B_{(k+1)\Delta_s} - B_{k\Delta_s}\}_{k=0}^{\Delta_s^{-1}-1}$ up-to time 1.
 4. **Output:** $(u_l, u_{l-1}) = ((x_{1,l}, v_{1,l}), (x_{1,l-1}, v_{1,l-1}))$ as generated in step 3..
-

Algorithm 5 : Sampling from coupled kernel $\tilde{Q}_{l,l-1}$.

1. **Input:** $l \in \{l_* + 1, l_* + 2, \dots\}$, and for $s \in \{l, l-1\}$,
 $(u_s, \tilde{u}_s) = ((x_{0,s}, v_{0,s}), (\tilde{x}_{0,s}, \tilde{v}_{0,s}))$.

2. Sample $(\Gamma_{0,l}, \dots, \Gamma_{\Delta_l^{-1}-1,l})$ and $(B_{\Delta_l}, B_{2\Delta_l} - B_{\Delta_l}, \dots, B_1 - B_{1-\Delta_l})$.
Concatenate to obtain $(\Gamma_{0,l-1}, \dots, \Gamma_{\Delta_{l-1}^{-1}-1,l-1})$
and $(B_{\Delta_{l-1}}, B_{2\Delta_{l-1}} - B_{\Delta_{l-1}}, \dots, B_1 - B_{1-\Delta_{l-1}})$.
 3. For $s \in \{l, l-1\}$: run the recursion (2.4)-(2.5) with $\{\Gamma_{k,s}\}_{k=0}^{\Delta_s^{-1}-1}$
and $\{B_{(k+1)\Delta_s} - B_{k\Delta_s}\}_{k=0}^{\Delta_s^{-1}-1}$ up-to time 1.
 4. **Output:** for $s \in \{l, l-1\}$, $(u_s, \tilde{u}_s) = ((x_{1,s}, v_{1,s}), (\tilde{x}_{1,s}, \tilde{v}_{1,s}))$ as generated in step 3.
-

Algorithm 6 : Sampling from coupled kernel $\check{P}_{l,l-1}$.

1. **Input:** $l \in \{l_* + 1, l_* + 2, \dots\}$, and for $s \in \{l, l-1\}$,
 $(u_s, \tilde{u}_s) = ((x_{0,s}, v_{0,s}), (\tilde{x}_{0,s}, \tilde{v}_{0,s}))$.
 2. Sample $(\Gamma_{0,l}, \dots, \Gamma_{\Delta_l^{-1}-2,l})$ and $(B_{\Delta_l}, B_{2\Delta_l} - B_{\Delta_l}, \dots, B_{1-\Delta_l} - B_{1-2\Delta_l})$.
Concatenate to obtain $(B_{\Delta_{l-1}}, B_{2\Delta_{l-1}} - B_{\Delta_{l-1}}, \dots, B_{1-\Delta_{l-1}} - B_{1-2\Delta_{l-1}})$
and $(\Gamma_{0,l-1}, \dots, \Gamma_{\Delta_{l-1}^{-1}-2,l-1})$.
 3. For $s \in \{l, l-1\}$: run the recursion (2.4)-(2.5) with $\{\Gamma_{k,s}\}_{k=0}^{\Delta_s^{-1}-1}$
and $\{B_{(k+1)\Delta_s} - B_{k\Delta_s}\}_{k=0}^{\Delta_s^{-1}-1}$ up-to time $1 - \Delta_s$.
 4. Sample

$$\left((X_{1,l}, V_{1,l}), (\tilde{X}_{1,l}, \tilde{V}_{1,l}) \right), \left((X_{1,l-1}, V_{1,l-1}), (\tilde{X}_{1,l-1}, \tilde{V}_{1,l-1}) \right) \Big| \left((x_{1-\Delta_l,l}, v_{1-\Delta_l,l}), (\tilde{x}_{1-\Delta_l,l}, \tilde{v}_{1-\Delta_l,l}), \right. \\ \left. (x_{1-\Delta_{l-1},l-1}, v_{1-\Delta_{l-1},l-1}), (\tilde{x}_{1-\Delta_{l-1},l-1}, \tilde{v}_{1-\Delta_{l-1},l-1}) \right)$$
- from the synchronous pairwise reflection maximal coupling [4].
5. **Output:** for $s \in \{l, l-1\}$, $(u_s, \tilde{u}_s) = ((x_{1,s}, v_{1,s}), (\tilde{x}_{1,s}, \tilde{v}_{1,s}))$ as generated in step 4.
-

References

- [1] BENVENISTE, A., METIVIER, N. & PRIOURET, P. (1990). *Adaptive Algorithms and Stochastic Approximations. Applications of Mathematics*. Springer: New York.
- [2] BESKOS, A., PAPASPILIOPOULOS, O., ROBERTS, G. O. & FEARNHEAD, P. (2006) Exact and computationally efficient likelihood-based estimation for discretely observed diffusion processes. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, **68**, pp. 333–382.
- [3] BESKOS, A. & ROBERTS, G. (2005). Exact simulation of diffusions, *Ann. Appl. Probab.*, **15**, 422–2444.

- [4] BOU-RABEE, N., EBERLE, A., & ZIMMER, R. (2020). Coupling and convergence for Hamiltonian Monte Carlo. *Ann. Appl. Probab.*, **30**, 1209–1250.
- [5] CHADA, N. K., FRANKS, J., JASRA, A., LAW, K. J. H. & VIHOLA, M. (2021) Unbiased estimation of discretely observed hidden Markov models. *SIAM/ASA JUQ*, 9(2), 763-787.
- [6] CHADA, N. K., JASRA, A. & YU, F. (2022). Unbiased estimation of the Hessian for partially observed diffusions. *Proceedings of the Royal Society A.*, **478**, 20210710.
- [7] CHENG, X., CHATTERJI, N. S. BARTLETT, P. L., & JORDAN, M. I. (2018). Underdamped Langevin MCMC: A non-asymptotic analysis. *Proceedings of Machine Learning Research*, **75**, 1–24.
- [8] CHIZAT, L. & BACH F. (2018). On the global convergence of gradient descent for overparameterized models using optimal transport. *Advances in neural information processing systems*, 3040–3050.
- [9] DALALYAN, A. S. (2017). Theoretical guarantees for approximate sampling from smooth and logconcave densities. *J. R. Statist. Soc. Ser. B*, **79** 651–676.
- [10] DUNCAN, A., NUSKEN, N. & PAVLIOTIS, G. A. (2017). Using perturbed underdamped Langevin dynamics to efficiently sample from probability distributions. *J Stat Phys.*, **169**, 1098–1131.
- [11] DURMUS, A., ENFROY, A. MOULINES, E., & STOLTZ, G. (2021). Uniform minorization condition and convergence bounds for discretizations of kinetic Langevin dynamics. arXiv preprint.
- [12] DURMUS, A. & MOULINES, E. (2016) Sampling from strongly log-concave distributions with the Unadjusted Langevin Algorithm. arxiv preprint.
- [13] FRENKEL, D. & SMIT, B. (2002). *Understanding Molecular Simulation, From Algorithms to Applications*. Academic Press: New York.
- [14] GAO, X., GURBUZBALABAN, M. & LINGJIONG, Z. (2020). Breaking reversibility accelerates langevin dynamics for non-convex optimization. *NIPS'20: Proceedings of the 34th International Conference on Neural Information Processing Systems*, **498**, 17850–17862.
- [15] GINZBURG, V.L. & LANDAU, L.D. (1950). *Zh. Eksp. Teor. Fiz.* **20**, 1064. English translation in: L. D. Landau, *Collected papers* (Oxford: Pergamon Press, 1965) p. 546
- [16] GLYNN, P. W. & RHEE, C. H. (2014). Exact estimation for Markov chain equilibrium expectations. *J. Appl. Probab.*, **51**, 377–389.
- [17] HENG, J., HOUSSINEAU, J. & JASRA, A. (2021). On unbiased score estimation for partially observed diffusions. arXiv preprint.
- [18] HENG, J. & JACOB, P. (2019). Unbiased Hamiltonian Monte Carlo with couplings. *Biometrika.*,106(2):287–302.

- [19] HENG, J., JASRA, A. LAW, K. J. H. & TARAKANOV, A. (2023). On unbiased estimation for discretized models. *SIAM/ASA JUQ* (to appear).
- [20] HOHENBERG, P.C. & KREKHOV, A.P. (2015). An introduction to the Ginzburg-Landau theory of phase transitions and nonequilibrium patterns. *Physics Reports*, **572**, 1–42.
- [21] JACOB, P., O’ LEARY, J. & ATCHADÉ, Y. (2020). Unbiased Markov chain Monte Carlo with couplings (with discussion). *J. R. Statist. Soc. Ser. B*, **82**, 543–600.
- [22] JASRA, A., LAW, K. J. H. & LU, D. (2021). Unbiased estimation of the gradient of the log-likelihood in inverse problems. *Stat. Comp.* **31**, 21.
- [23] JASRA, A., LAW, K. J. H. & YU, F. (2022). Unbiased filtering of a class of partially observed diffusions. *Adv. Appl. Probab.* **54**, 661-687.
- [24] JASRA, A., LAW, K. J. H. & XU, Y. (2021). Markov chain Simulation for Multilevel Monte Carlo. *Found. Data Sci.* **3**, 27-47.
- [25] JELIC, V. & MARSIGLIO, F. (2012) The double-well potential in quantum mechanics: a simple, numerically exact formulation. *Eur. J. Phys.* **33** 1651.
- [26] KAZEYKINA, A., REN, Z. TAN, X. & YANG, J. (2020) Ergodicity of the underdamped mean-field Langevin dynamics. arXiv preprint.
- [27] LIVINGSTONE, S., FAULKNER, M. F. & ROBERTS, G. O. (2019). Kinetic energy choice in Hamiltonian/hybrid Monte Carlo. *Biometrika*, **106**, 2, 303–319.
- [28] MEI, S., A. MONTANARI, A, & P.-M. NGUYEN P.-M. (2018), A mean field view of the landscape of two-layer neural networks, *Proceedings of the National Academy of Sciences*, **115**, 7665–7671.
- [29] MCLEISH, D. (2011). A general method for debiasing a Monte Carlo estimator. *Monte Carlo Meth. Appl.*, **17**, 301–315.
- [30] MIDDLETON, L., DELIGIANNIDIS, G., DOUCET, A. & JACOB, P. (2019). Unbiased smoothing using particle independent Metropolis-Hastings. In *Proc. Mach. Learn. Res.*, **89**.
- [31] NEAL, R. M. (2011). MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*. CRC Press.
- [32] PAVLIOTIS, G. A. (2014). *Stochastic Processes and Applications*. Springer.
- [33] PROPP, J. G., & WILSON, D. B. (1996). Exact sampling with coupled Markov chains and applications to statistical mechanics. *Rand. Struct. Algs.*, **9**, 223–252.
- [34] RHEE, C. H. (2013). Unbiased estimation with biased samples. Ph.D. thesis, Stanford University.
- [35] RHEE, C. H. & GLYNN, P. (2015). Unbiased estimation with square root convergence for SDE models. *Op. Res.*, **63**, 1026–1043.

- [36] ROBERT, C. & CASELLA, G. (2004). *Monte Carlo Statistical Methods*. Springer: New York.
- [37] ROBERTS, G. O. & STRAMER, O. (2003) Langevin diffusions and Metropolis-Hastings algorithms. *Methodol. Comput. Appl. Probab.*, 4(4):337–357.
- [38] RUZAYQAT, H., BESKOS, A., CRISAN, D., JASRA, A. & KANTAS, N. (2023). Unbiased Estimation using a Class of Diffusion Processes. *J. Comp. Phys.*, **472**, 111643
- [39] RUZAYQAT, H. & JASRA A. (2022). Unbiased Parameter Inference for a Class of Partially Observed Lévy-Process Models. *Found. Data Sci.*, **4**(2), 299–322.
- [40] RUZAYQAT, H. & JASRA A. (2020) Unbiased estimation of the solution to Zakai’s equation. *Monte Carlo Meth. Appl.*, **26**(2), 113–129.
- [41] SIMSEKLI, U., LINGJIONG, Z., WHY TEH, Y. & GURBUZBALABAN, M. (2020) Fractional underdamped Langevin dynamics: retargeting SGD with momentum under heavy-tailed gradient noise. *Proceedings of the 37th International Conference on Machine Learning*, **832**, 8970–8980.
- [42] THORISSON, H. (2000). *Coupling, Stationarity and Regeneration*. Springer: New York.
- [43] VAN DEN BOOM, W., JASRA, A., DE IORIO, M., & BESKOS, A. (2022). Unbiased approximation of posteriors via coupled particle Markov chain Monte Carlo. *Stat. Comp.* **32**, article 36.
- [44] VIHOLA, M. (2018). Unbiased estimators and multilevel Monte Carlo. *Operations Research*, **66**(2), 448–462.
- [45] WANG, G., O’LEARY, J. & JACOB, P. (2021). Maximal Couplings of the Metropolis-Hastings Algorithm. *International Conference on Artificial Intelligence and Statistics*, 1225–1233.