# Unbiased Multilevel Monte Carlo methods for intractable distributions: MLMC meets MCMC

Guanyang Wang[*] and Tianze Wang[†]

March 23, 2022

**Abstract**

Constructing unbiased estimators from MCMC outputs has recently increased much attention in statistics and machine learning communities. However, the existing unbiased MCMC framework only works when the quantity of interest is an expectation of certain intractable probability distribution. In this paper, we propose unbiased estimators for functions of expectations. Our idea is based on the combination of the unbiased MCMC and MLMC methods. In contrast to traditional sequential methods, our estimator can be implemented on parallel processors independently. We prove our estimator has a finite variance, a finite computational complexity, and achieves $\varepsilon$-accuracy with $O(1/\varepsilon^2)$ computational cost under mild conditions. We also illustrate our estimator on several numerical examples.

## 1  Introduction

Monte Carlo methods provide unbiased estimators for the expected value of a distribution. In practice, however, the distribution may be infeasible to sample from, and the quantity of interest may not be an expected value. Generally speaking, most inference problems can be viewed as estimating a quantity of the form $\mathcal{T}(\pi)$, where $\pi$ stands for one or a family of distributions, and $\mathcal{T}$ is a scalar or vector-valued functional of $\pi$. We start with the following motivating examples to gain further insights.

---

[*]Department of Statistics, Rutgers University, New Brunswick, USA. Email: guanyang.wang@rutgers.edu

[†]Department of Statistics, Rutgers University, New Brunswick, USA. Email: tw522@scarletmail.rutgers.edu

Alphabetical authorship.

**Example 1** (Estimating integration). *Let $\pi$ be a probability distribution and $f$ a $\pi$-integrable function. The problem of estimating $\int f(x)\pi(dx)$ can be viewed as estimating $\mathcal{T}(\pi)$ where $\mathcal{T}$ is the integral operator:*

$$\mathcal{T}(\pi) := \int f(x)\pi(dx).$$

**Example 2** (Nested Monte Carlo). *Let $\pi$ be a probability distribution, and suppose the quantity of our interest has the form $\mathcal{T}(\pi) := \mathbb{E}_\pi[\lambda]$, where $\lambda$ is itself intractable. The intractable function $\lambda$ may take the form $\lambda(x) := f(x, \gamma(x))$, where $\gamma(x) = \mathbb{E}_{y\sim p(y|x)}[\phi(x,y)]$ is an expectation. One concrete example is the two-stage optimal stopping problem where $\gamma(x) = \max\{x, \mathbb{E}[y|x]\}$. Estimating nested expectation is a common problem in statistics [Giles and Goda, 2019], machine learning [Rainforth et al., 2018], and operation research [Belomestny et al., 2015, Zhou et al., 2021]. It is known that the standard Monte Carlo estimator generally has systematic bias and suboptimal convergence rate.*

**Example 3** (Simulating ratios of normalizing constants). *Let*

$$\pi_1(x) = \frac{f_1(x)}{Z_1}$$

*and*

$$\pi_2(x) = \frac{f_2(x)}{Z_2}$$

*be two probability densities with common support. We assume $f_1$ and $f_2$ can be easily evaluated, but the normalizing constants $Z_1$ and $Z_2$ are computationally intractable. Consider the task of estimating the ratio of normalizing constants, i.e., $Z_1/Z_2$, standard calculation yields:*

$$\frac{Z_1}{Z_2} = \frac{\mathbb{E}_{\pi_2}[f_1]}{\mathbb{E}_{\pi_1}[f_2]}. \tag{1}$$

*The problem can be viewed as estimating $\mathcal{T}(\pi)$ by choosing $\pi$ as the product measure $\pi_1 \times \pi_2$, and*

$$\mathcal{T}(\pi) := \frac{\mathbb{E}_{\pi_2}[f_1]}{\mathbb{E}_{\pi_1}[f_2]}.$$

*The problem finds many statistical and physics applications including hypothesis testing, Bayesian inference, and estimating free energy differences. We refer the readers to Meng and Wong [1996] for other interesting applications.*

**Example 4** (Quantile estimation). *Let $\pi$ be a probability distribution with cumulative distribution function $F_\pi$ and $q$ a constant in $(0,1)$. Estimating the $q$-th quantile of $\pi$ can be formulated as estimating $\mathcal{T}(\pi)$ where*

$$\mathcal{T}(\pi) := \inf_v\{F_\pi(v) \geq q\}.$$

*Quantile estimation problem has widespread applications in statistics, machine learning, economics, and other fields. We refer the readers to Koenker and Hallock [2001], Takeuchi et al. [2006], Romano et al. [2019] for more discussions, and Doss et al. [2014] for an MCMC-based quantile estimation method.*

In all the aforementioned examples, the distribution $\pi$ can be intractable. In some cases, such as Example 1 and 2, the quantity of interest is an expectation under $\pi$, although the function inside the expectation may or may not be intractable. In other cases, including Example 3 and 4, $\mathcal{T}$ is a functional of $\pi$, but not an expectation.

Throughout this paper, we focus on the unbiased estimation of $\mathcal{T}(\pi)$ assuming one can only get access to outputs from some MCMC algorithm that leaves $\pi$ as stationary distribution. There are many reasons that we are interested in unbiased estimators. Firstly, classical MCMC estimators are accurate until the number of iterations reaches infinity. Estimators based on empirical averages of MCMC outputs are generally biased for any fixed number of iterations. The presence of bias precludes the direct use of modern parallel computing architectures, indicated by an increasing number of processors but a limited computational budget per processor. Averaging independent chains in parallel will only decrease the variance, but the bias remains the same. In contrast, unbiased estimators can be computed on different devices in parallel without communication. On top of parallel computing, the confidence intervals can be easily constructed using unbiased estimators from Monte Carlo outputs, leading to better uncertainty quantification in settings where the variance is difficult to estimate. Moreover, unbiased estimators are often used as a component in several classes of exact inference methods such as the pseudo-marginal Metropolis–Hastings algorithms [Andrieu and Roberts, 2009]. In addition, unbiasedness is appealing in various practical applications, including Rischard et al. [2018], Chen et al. [2018], Tadić and Doucet [2017] and the references therein.

Without further assumption on $\mathcal{T}$ and $\pi$, unbiased estimation of $\mathcal{T}(\pi)$ is generally a challenging problem. Computational challenges appear in both components of the pair $(\mathcal{T}, \pi)$. Fortunately, recent works provide promising solutions when one component of the above $(\mathcal{T}, \pi)$ pair is easy while the other is relatively difficult. We briefly discuss the following two cases separately:

- (Case 1: Easy $\mathcal{T}$, difficult $\pi$): When $\mathcal{T}$ is an integral operator with respect to some tractable function $f$, but $\pi$ is infeasible to sample from, i.e., $\mathcal{T}(\pi) := \mathbb{E}_\pi[f]$ for some intractable $\pi$. The problem is considered by Jacob, O'Leary and Atchadé (JOA henceforth) [Jacob et al., 2020]. The JOA estimator, which follows the idea of Glynn and Rhee [2014], solves this problem via couplings of Markov chains. The unbiased MCMC framework has recently raised much attention. It has been applied in convergence diagnostics [Biswas et al., 2019, Biswas and Mackey, 2021], gradient estimation [Ruiz et al., 2020], asymptotic variance estimation [Douc et al., 2021], and so on.
- (Case 2: Easy $\pi$, difficult $\mathcal{T}$): When $\pi$ can be sampled perfectly, but $\mathcal{T}(\pi) := g(\mathbb{E}_\pi[f])$ is a function of the expectation or the nested expectation $\mathbb{E}_{x\sim\pi}[\gamma(x, \mathbb{E}_{y|x}[\phi(y)])]$, the state of the art debiasing technique is the unbiased Multilevel Monte Carlo (MLMC) method developed by McLeish, Glynn, Rhee, and Blanchet [Blanchet et al., 2015, Rhee and Glynn, 2015, McLeish, 2011] which is a randomized version of the celebrated MLMC methods pioneered by Heinrich and Giles [Heinrich, 2001, Giles, 2008, 2015]. Unbiased MLMC methods have also found many applications including gradient estimation [Shi and Cornish, 2021], optimal stopping [Zhou et al., 2021], robust optimization [Levy et al., 2020].

In summary, the unbiased MCMC method assumes easy $\mathcal{T}$ (an integral operator) but difficult $\pi$, and the unbiased MLMC method assumes easy $\pi$ (perfectly simulable) but

difficult $\mathcal{T}$. Both assumptions can be violated in many practical applications, such as Example 2, 3 and 4. Despite immense progress that has been made in the debiasing techniques, there is no systematic way of constructing unbiased estimators of $\mathcal{T}(\pi)$ beyond special cases.

In this article, we present a step towards designing unbiased estimators of $\mathcal{T}(\pi)$ for general $(\mathcal{T}, \pi)$ pair by combining the ideas of the unbiased MCMC and MLMC methods. We propose generic unbiased estimators of functional of expectations, i.e.,

$$T(\pi) = g(m(\pi)) := g(\mathbb{E}_{\pi}[f(X)])$$

where $\pi$ is a $d$-dimensional probability measure, $f : \mathbb{R}^d \to \mathbb{R}^m$ is a deterministic map, and $g : \mathbb{R}^m \to \mathbb{R}$ is a deterministic function [1]. Other technical assumptions, including the smoothness of $g$ and the existence of moments, will be made clear in the subsequent sections. This estimator can be naturally extended to unbiasedly estimate nested expectations, as will be discussed in Section 3.2 in details. The unbiased estimator is easily parallelizable. It has both finite variance and computational cost for a general class of problems, which implies a 'square root convergence rate' that matches the oracle rate given by the Central Limit Theorem. Moreover, some technical assumptions on $g$ relax the standard 'linear growth' assumption in Blanchet and Glynn [2015] and Blanchet et al. [2019], which may be of independent interest.

Our method naturally connects the unbiased MCMC with the MLMC method, which are developed in their own communities. Unbiased MCMC is an emerging area for its potential of parallelization. The methodology in Jacob et al. [2020] has been extended to different MCMC algorithms including the Hamiltonian Monte Carlo [Heng and Jacob, 2019] and the pseudo-marginal MCMC [Middleton et al., 2020]. In contrast, the MLMC method (both the non-randomized and randomized version) is shown to be successful in estimating the expectation of SDE solutions [Giles, 2008, Rhee and Glynn, 2015], option pricing [Belomestny et al., 2015, Zhou et al., 2021], and inverse problems [Beskos et al., 2017, Jasra et al., 2018]. It seems that the unbiased MCMC method is mostly developed and applied in statistics and machine learning, while the MLMC method is often used in applied mathematics, operation research, and computational finance. We hope this work will serve as a bridge for these communities, and invite researchers from broader areas to develop these methods together.

The rest of this paper is organized as follows. Section 1.1 introduces the notations. In Section 2 we describe the high-level idea behind our method without diving into details. The connections between unbiased MCMC and MLMC methods will also be clear in this section. In Section 3 we formally propose our estimator, state the assumptions, and prove its theoretical properties. In Section 4 we implement our method on several examples to study its empirical performance. We conclude this paper in Section 5. Technical details such as proofs and implementation details are deferred to the Appendix.

---

[1]For simplicity, we only consider scalar-valued $g$ in this paper, though our method can be naturally generalized to vector-valued functions.

## 1.1 Notations

Throughout this article, we preserve the notation $g$ to denote a function from its domain $\mathcal{D} \subset \mathbb{R}^m$ to $\mathbb{R}$ of our interest. The domain $\mathcal{D}$ plays an important role in both algorithm design and theoretical analysis, which will be addressed in Section 3. We write $\pi$ as a $d$-dimensional probability measure, and $\pi_1, \cdots, \pi_d$ for its marginal distributions. We denote by $m_f(\pi) := \mathbb{E}_\pi[f(X)]$ the expected value/vector of $f$ under $\pi$, and write it as $m(\pi)$ when it is unlikely to cause confusion. The $L^p$ norm of a vector $v \in \mathbb{R}^d$ is written as $\|v\|_p := \left( \sum_{i=1}^d v_i^p \right)^{1/p}$. For the $L^2$ norm, we simply write $\|v\| := \|v\|_2 = \sqrt{\sum_{i=1}^d v_i^2}$. The geometric distribution with success probability $r$ is denoted by $\mathsf{Geo}(r)$, and $p_n = p_n(r) := \mathbb{P}(\mathsf{Geo}(r) = n) = (1-r)^{n-1} r$. The uniform distribution on $[0,1]$ is denoted by $\mathsf{U}[0,1]$. The multivariate normal distribution with mean $\mu \in \mathbb{R}^d$ and covariance matrix $\Sigma$ is denoted by $\mathsf{N}(\mu, \Sigma)$. The binomial distribution with $N$ trials and parameter $p$ is denoted by $\mathsf{Binom}(N, p)$. The Poisson distribution with distribution $\lambda$ is denoted by $\mathsf{Poi}(\lambda)$. We adopt the convention that $\sum_{i=m}^n a_i = 0$ if $m > n$. Given a set $A \subset \mathbb{R}^d$, we denote by $A^\circ$ all the interior points of $A$. For a differentiable function $h : \mathbb{R}^d \to \mathbb{R}$, we denote by $Dh := (\frac{\partial h}{\partial x_1}, \frac{\partial h}{\partial x_2}, \cdots, \frac{\partial h}{\partial x_d})$ the gradient of $h$. Given two probability measures $\mu$ and $\nu$ on the same probability space, we write their total variation (TV) distance as $\|\mu - \nu\|_{\mathsf{TV}} := \sup_A |\mu(A) - \nu(A)|$.

## 2 A Simple Identity: Unbiased MCMC meets MLMC

Consider the task of designing unbiased estimators of $g(m(\pi)) = g\big(\mathbb{E}_\pi[f(X)]\big)$. The problem is extensively studied in the literature when one can draw independent and identically distributed ($i.i.d.$) samples from $\pi$. Unbiased estimators are known to exist, or not exist under different contexts [Keane and O'Brien, 1994, Jacob and Thiery, 2015]. Different debiasing techniques including Bernstein polynomials [Nacu and Peres, 2005], Taylor polynomials [Blanchet et al., 2015], and the MLMC method [Blanchet and Glynn, 2015, Vihola, 2018] have been proposed and analyzed. The core idea in all these techniques is to write $g(m(\pi))$ as an infinite series $\sum_{k=1}^\infty a_k$. One can then choose a random level $k$ with probability $p_k$ and construct the importance sampling-type estimator $\hat{a}_k / p_k$. Suppose each $\hat{a}_k$ is unbiased for $a_k$, then $\hat{a}_k / p_k$ is generally unbiased for $\sum_{k=1}^\infty a_k$. Among the existing methods, the unbiased MLMC framework seems to work with the greatest generality, as it does not require the knowledge on the derivatives of $g$.

When $\pi$ is infeasible to sample from, our first observation is based on the following simple identity:

$$g(m(\pi)) = g(m(\tilde{\pi})) \tag{2}$$

for every $\tilde{\pi}$ which satisfies $m(\tilde{\pi}) = m(\pi)$. Formula (2) is straightforward mathematically, but it turns out the right hand side of (2) is more tractable than the left hand side. To be more precise, one main difficulty in estimating $g(m(\pi))$ arises from the difficulty in sampling $\pi$. However, one observation is the quantity $g(m(\pi))$ essentially depends only an expectation under $\pi$ but not directly on $\pi$ itself. Therefore, formula (2) suggests that one can sample unbiased estimators of $m(\pi)$ instead of sampling from $\pi$ directly.

If we are able to *i.i.d.* sample unbiased estimators of $m(\pi)$, which means we have *i.i.d.* samples from some distribution $\tilde{\pi}$ with $m(\tilde{\pi}) = m(\pi)$, then it suffices to estimate $g(m(\tilde{\pi}))$ unbiasedly. The difficulty is now reduced to estimating a functional of expectation, and the existing unbiased MLMC methods (and all the debiasing tricks mentioned above) can be directly applied.

After observing (2), it suffices to construct unbiased estimators of $m(\pi)$ provided that $\pi$ cannot be directly simulated. The unbiased MCMC framework provides us with natural solutions. Suppose a Markov chain with transition kernel $P$ that targets $\pi$ as stationary distribution. It is often possible to construct a pair of coupled Markov chains $(Y, Z) = (Y_t, Z_t)_{t=1}^{\infty}$ that both evolve according to $P$. By design, if the pair $(Y_t, Z_{t-1})$ meets at some random time $\tau$ and stay together after meeting, then the Jacob-O'Leary-Atchadé (JOA) estimator is unbiased for $m(\pi)$. Therefore, we can unbiasedly estimate $g(m(\pi))$ using the following two-step simulation strategy. The overall workflow is described in the Figure 1 below. The unbiased MCMC algorithm is used here as a generator for random variables with expectation $m(\pi)$. Therefore we can use the outputs of the unbiased MCMC algorithm as inputs for the unbiased MLMC approach, and eventually construct an unbiased estimator of $g(m(\pi))$.
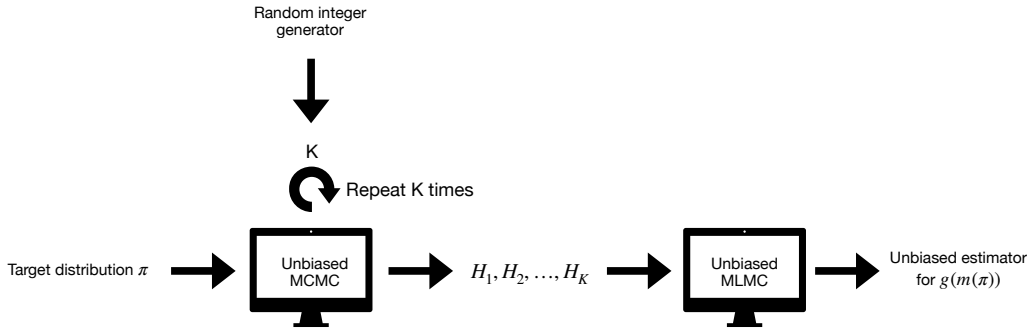


Figure 1: The workflow for constructing an unbiased estimator of $g(m(\pi))$.

# 3 Unbiased estimators for functions of expectation

In this section, we discuss our estimator for $g(m(\pi))$ from MCMC outputs in details. We start with a brief review on the JOA estimator of $m(\pi)$ in Section 3.1.1. Our general framework is described in Section 3.1.2. A family of much simplified estimators is given in Section 3.1.3 when $g$ admits additional structures. In Section 3.2 we discuss the generalization of our approach to unbiasedly estimate nest expectations. In Section 3.3 we discuss the problem regarding the domain of $g$. We provide a transformation that avoids the domain problem for a general class of functions. In Section 3.4 we give theoretical justifications of our method.

## 3.1 Constructing the unbiased estimator

### 3.1.1 The Jacob-O'Leary-Atchadé (JOA) estimator of $m(\pi)$

Let $\Omega$ be a Polish space equipped with the standard Borel $\sigma$-algebra $\mathcal{F}$. Let $P : \Omega \times \mathcal{F} \to [0, 1]$ be the Markov transition kernel that leaves $\pi$ as stationary distribution. The Jacob-O'Leary-Atchadé (JOA) estimator uses a coupled pair of Markov chains that both has transition kernel $P$. Formally, the coupled pair $(Y, Z) = (Y_t, Z_t)_{t=1}^{\infty}$ can be viewed as a Markov chain on the produce space $\Omega \times \Omega$. The transition kernel $\bar{P}$, which is also called the coupling of $(Y, Z)$, satisfies

$$\bar{P}((x,y), A \times \Omega) = P(x, A), \bar{P}((x,y), \Omega \times B) = P(y, B)$$

for every $x, y \in \Omega$ and $A, B \in \mathcal{F}$. The coupled chain starts with $Y_0 \sim \pi_0, Y_1 \sim P(Y_0, \cdot)$ and $Z_0 \sim \pi_0$ independently. Then at each step $t \geq 2$, one samples $(Y_t, Z_{t-1}) \sim \bar{P}((Y_{t-1}, Z_{t-2}), \cdot)$. Suppose the coupling $\bar{P}$ is 'faithful' [Rosenthal, 1997], meaning that there is a random but almost surely finite time $\tau$ such that $Y_\tau = Z_{\tau-1}$, and $Y_t = Z_{t-1}$ for every $t \geq \tau$. Then for every $k$,

$$H_k(Y, Z) := f(Y_k) + \sum_{i=k+1}^{\tau-1} (f(Y_i) - f(Z_{i-1})) \tag{3}$$

is an unbiased estimator. For any fixed integer $m \geq k$, the following 'time-averaged' estimator $H_{k:m}(Y, Z) := (m - k + 1)^{-1} \sum_{l=k}^{m} H_l(Y, Z)$ is the average of $m - k + 1$ unbiased estimators, which retains unbiasedness and reduces the variance. The unbiasedness of $H_k(Y, Z)$ is justified by the following informal calculation in Jacob et al. [2020]:

$$m(\pi) = \lim_{n \to \infty} \mathbb{E}[f(Y_n)] = \mathbb{E}[f(Y_k)] + \sum_{n=k+1}^{\infty} (\mathbb{E}[f(Y_n)] - \mathbb{E}[f(Y_{n-1})])$$

$$= \mathbb{E}[f(Y_k)] + \sum_{n=k+1}^{\infty} \mathbb{E}[f(Z_n) - f(Y_{n-1})]$$

$$= \mathbb{E}[f(Y_k)] + \sum_{n=k+1}^{\tau-1} \mathbb{E}[f(Z_n) - f(Y_{n-1})] = \mathbb{E}[H_k(Y, Z)].$$

The rigorous proof requires assumptions on the target $\pi$, and the distribution of $\tau$. We refer the readers to Jacob et al. [2020], Middleton et al. [2020] and our appendix for

more details. Theoretical and empirical investigation of these are provided in O'Leary and Wang [2021], Wang et al. [2021]. More sophisticated unbiased estimators using $L$-lag coupled chains are discussed in Biswas et al. [2019], but the main idea remains the same.

### 3.1.2 Unbiased estimator of $g(m(\pi))$

Suppose we get access to a routine $\mathcal{S}$ such as the JOA estimator in Section 3.1.1 which outputs unbiased estimators of $m(\pi)$. The estimator of $g(m(\pi))$ can then be constructed by the randomized MLMC method. Let $H_1, H_2, \cdots, H_{2m}$ be a sequence of $i.i.d.$ random variables. We let

$$S_H(2m) := \sum_{k=1}^{2m} H_i \tag{4}$$

be the summation of all the $2m$ terms, and let

$$S_H^{\mathsf{O}}(m) := \sum_{k=1}^{m} H_{2k-1} \qquad S_H^{\mathsf{E}}(m) := \sum_{k=1}^{m} H_{2k} \tag{5}$$

be the summation of all the odd and even terms, respectively. The estimator is described by Algorithm 1.

---

**Algorithm 1** Unbiased Multilevel Monte-Carlo estimator

---

**Input:**
- A subroutine $\mathcal{S}$ for generating unbiased estimators of $m(\pi)$
- A function $g : \mathcal{D} \to \mathbb{R}$
- The parameter $p$ for geometric distribution

**Output:** Unbiased estimator of $g(m(\pi))$
  1. Sample $N$ from the geometric distribution $\mathsf{Geo}(p)$
  2. Call $\mathcal{S}$ for $2^N$ times and label the outputs by $H_1, ..., H_{2^N}$
  3. Calculate the quantities $S_H(2^N)$, $S_H^{\mathsf{O}}(2^{N-1})$ and $S_H^{\mathsf{E}}(2^{N-1})$ by (4),(5)
  4. Calculate $\Delta_N = g\left(S_H(2^N)/2^N\right) - \frac{1}{2}\left(g\left(S_H^{\mathsf{O}}(2^{N-1})/2^{N-1}\right) + g\left(S_H^{\mathsf{E}}(2^{N-1})/2^{N-1}\right)\right)$

**Return:** $W = \Delta_N/p_N + g(H_1)$.

---

Again, the following informal calculation justifies the unbiasedness of $W$.

$$
\begin{aligned}
\mathbb{E}\left[W\right] &= \mathbb{E}\left[g(H_1)\right] + \mathbb{E}\left[\Delta_N/p_N\right] \\
&= \mathbb{E}\left[g(H_1)\right] + \mathbb{E}[\mathbb{E}\left[\Delta_N/p_N \mid N\right]] \\
&= \mathbb{E}\left[g(H_1)\right] + \sum_{n=1}^{\infty} \mathbb{E}[\Delta_n] \\
&= \mathbb{E}\left[g(H_1)\right] + \sum_{n=1}^{\infty} (\mathbb{E}[g(S_H(2^n)/2^n)] - \mathbb{E}[g(S_H(2^{n-1})/2^{n-1})]) \\
&= \lim_{n\to\infty} \mathbb{E}[g(S_H(2^n)/2^n)] = g(m(\tilde{\pi})) = g(m(\pi)),
\end{aligned}
$$

where $\tilde{\pi}$ is the distribution of each $H_i$.

Algorithm 1 differs from the original unbiased MLMC design Blanchet and Glynn [2015] as it relaxes the assumption '*i.i.d.* samples from $\pi$' by 'unbiased estimator of $m(\pi)$'. It also provides practical methods for finding the estimator via the JOA estimator Jacob et al. [2020] described in Section 3.1.1. It is worth mentioning that, any unbiased estimator of $m(\pi)$, including but not limited to the JOA estimator (see also Glynn and Rhee [2014], Agapiou et al. [2018]), can be fed into Algorithm 1 as a subroutine. On the other hand, we find the JOA estimator is by far the most general framework for constructing unbiased estimators of $m(\pi)$ given intractable $\pi$. We will implicitly assume the subroutine $\mathcal{S}$ is the JOA estimator in the subsequent sections.

### 3.1.3 Unbiased estimator of polynomials and other special functions

Section 3.1.2 provides us a relatively general framework for unbiased estimators of $g(m(\pi))$. In some situations where the target function $g$ has certain nice properties, the unbiased estimators can be easily obtained without resorting to the unbiased MLMC framework. For example, if $g(x) = x^k$ is a univariate monomial function, one can simply call the unbiased MCMC algorithm $k$ times and obtain $H_1, \cdots, H_k$. The estimator $\prod_{l=1}^{k} H_l$ will then be unbiased for $m(\pi)^k$. The argument above can be naturally extended to the case where $m(\pi) \in \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}$ is a multivariate polynomial function. We use the multi-index $k = (k_1, \cdots, k_m)$ with $\sum_{i=1}^{m} k_i \leq n$ where $k_1, \ldots, k_m$ are non-negative integers, and $x^k = x_1^{k_1} x_2^{k_2} \cdots x_m^{k_m}$. Let $g(x) = \sum_{k \leq n} \alpha_k x^k$ denote a multivariate polynomial with degree at most $n$, the unbiased estimator of $g(m(\pi))$ can be constructed as follows. First, we call the unbiased MCMC subroutine $\mathcal{S}$ for $n$ times and label the outputs by $H_1, \cdots, H_n$, each is an independent unbiased estimator of $m(\pi)$. Then for each $k = (k_1, \cdots, k_m)$ we calculate the quantity

$$\hat{H}(k) = \prod_{l_1=1}^{k_1} H_{l_1,1} \prod_{l_2=k_1+1}^{k_1+k_2} H_{l_2,2} \cdots \prod_{l_m=k_1+\cdots+k_{m-1}+1}^{k_1+\cdots+k_m} H_{l_m,m},$$

where $H_{a,b}$ stands for the $b$-th coordinate of $H_a \in \mathbb{R}^d$. It is then clear from the independence of $H_1, \cdots, H_n$ that $\mathbb{E}[\hat{H}(k)] = m(\pi)^k$. Finally we output $\sum_k \alpha_k \hat{H}(k)$ which is unbiased for $g(m(\pi))$ by the linearity of expectation. It is different from Algorithm 1 as it requires a fixed number of calls for $\mathcal{S}$, and it does not require calculating the consecutive difference $\Delta_k$.

When $g : \mathbb{R} \to \mathbb{R}$ is a real analytic function on $\mathcal{D}$, i.e., $g(x) = \sum_{n=0}^{\infty} a_i (x - a)^n$ for some real number $a$. Suppose $\tilde{N}$ is a non-negative integer-valued random variable with $\mathbb{P}(\tilde{N} = k) = q_k$. The unbiased estimator for $g(m(\pi))$ can be constructed by first generating $\tilde{N}$, and then call the unbiased MCMC subroutine $\mathcal{S}$ for $\tilde{N}$ times with outputs $H_1, \cdots H_{\tilde{N}}$. The final estimator can be expressed by

$$\frac{a_{\tilde{N}}}{q_{\tilde{N}}} \cdot \frac{\prod_{j=1}^{\tilde{N}} (H_j - a)}{\tilde{N}!}$$

which follows from the idea in Blanchet et al. [2015]. In particular, when $g(x) = e^x$ and $\tilde{N}$ follows from the Poisson distribution, the estimator is known as the 'Poisson estimator' which is used in both physics and statistics, see Wagner [1987], Papaspiliopoulos [2009], Fearnhead et al. [2010] for more discussions. Albeit useful in many cases,

the power-series-type estimators generally have strong assumptions on the smoothness of the target function. It also requires the knowledge of all the higher-order derivatives of $g$, which is generally infeasible when $g$ is either a complicated univariate function, or a multivariate function.

Therefore, throughout this paper, we will mostly focus on using the unbiased MLMC framework for estimating $g(m(\pi))$ given its greatest generality. This subsection intends to remind our readers that there may be simper and easier choices when the target function $g$ behaves 'nice' enough.

## 3.2   Nested Expectations

Now we explain how our method can be extended to estimate the nested expectations. Recall that a nested expectation can be written as $\mathbb{E}_\pi[\lambda]$. The function $\lambda$ is assumed to be intractable with the form $\lambda(x) := f(x, \gamma(x))$, where $\gamma(x) = \mathbb{E}_{y \sim \pi(y|x)}[\phi(x, y)]$ is an expectation under the conditional distribution.

We first decompose the joint distribution $\pi(x, y)$ as the marginal distribution $\pi(x)$ times the conditional distribution of $\pi(y|x)$. Observe that when $x = x_0$ is fixed, then

$$\lambda(x_0) = f(x_0, \mathbb{E}_{y \sim \pi(y|x_0)}[\phi(x_0, y)])$$

is essentially a function of $\mathbb{E}_{\pi(y|x_0)}[\phi(x_0, y)]$. Therefore our unbiased estimator can be constructed as follows:

---
**Algorithm 2** Unbiased Multilevel Monte-Carlo estimator for nested expectation
---

1. Sample $x$ from $\pi(x)$
2. Given $x$, generate an unbiased estimator $\hat{f}(x)$ of $f(x, \mathbb{E}_{y \sim \pi(y|x)}[\phi(x, y)])$ using Algorithm 1

   **Return:**  $\hat{f}(x)$.

---

It can be checked directly that Algorithm 2 outputs an unbiased estimator of $\mathbb{E}_\pi[\lambda]$.

**Proposition 1.** *We have $\mathbb{E}[\hat{f}] = \mathbb{E}_\pi[\lambda]$.*

*Proof.* Using the law of iterated expectation, the expectation of $\hat{f}$ can be written as

$$\begin{aligned}
\mathbb{E}[\hat{f}] &= \mathbb{E}[\mathbb{E}[\hat{f}(x)|x]] \\
&= \int \mathbb{E}[\hat{f}(x)|x]\pi(dx) \\
&= \int f(x, \mathbb{E}_{y \sim \pi(y|x)}[\phi(x, y)])\pi(dx) \\
&= \mathbb{E}_\pi[\lambda].
\end{aligned}$$

$\square$

Algorithm 2 is useful when $\pi(x)$ can be directly sampled from, and $\pi(y|x)$ can be approximated sampled from some MCMC algorithms. This is satisfied in many practical

models, including the cut distributions that arises frequently in Bayesian modular inference [Plummer, 2015, Pompe and Jacob, 2021].

We conclude this subsection by presenting a typical example of the nested expectation, namely estimating the expected utility under partial information [Giles, 2018, Giles and Goda, 2019]. Other examples including the Bayesian experimental design and variational autoencoders are given in Rainforth et al. [2018], Hironaka and Goda [2021], Goda et al. [2022].

**Example 5** (The expected utility under partial information)**.** *Suppose we have an underlying two-stage process $(X, Y)$ with joint distribution $\pi(x, y) = \pi(x)\pi(y|x)$. Suppose we have $D$ possible strategies (for example, treatments), each with corresponding utility $f_d(x, y)$ for $d \in \{1, 2, \cdots, D\}$. Suppose we have to choose a strategy without seeing the values of $(X, Y)$, the optimal expected utility would be $\max_d \mathbb{E}[f_d(X, Y)]$. Similarly, after observing $(X, Y)$, the optimal strategy would be $d^*(X, Y) = \arg\max_d f_d(X, Y)$, with expected utility $\mathbb{E}[\max_d f_d(X, Y)]$. In the intermediate case, if one has observed $X$ but not $Y$, then the optimal strategy would maximize the conditional utility, i.e., $d^*(X) = \arg\max_d \mathbb{E}[f_d(X, Y)|X]$. The optimal strategy with partial information has expected utility $\mathbb{E}[\max_d \mathbb{E}[f_d(X, Y)|X]]$, which is a nested expectation.*

*The expected utility under partial information finds widespread applications in computational finance, especially in European option pricing Belomestny et al. [2015], Zhou et al. [2021]. Meanwhile, the difference between full and partial utility*

$$\mathbb{E}[\max_d f_d(X, Y)] - \mathbb{E}[\max_d \mathbb{E}[f_d(X, Y)|X]]$$

*quantifies the 'value' of the information in $Y$, which also has applications in the evaluation of Value-at-Risk (VaR) [Giles, 2018] and medical areas [Ades et al., 2004]. Existing literatures typically assume one can sample directly from the joint distribution $\pi(x, y)$, and regard the intractable $\pi(x, y)$ (such as posterior distribution) as an open question, see Section 5 of Giles and Goda [2019] for discussions.*

## 3.3  The domain problem and the $\delta$-transformation

There is an extra subtly in implementing Algorithm 1. Besides requiring $H$ to be an unbiased estimator of $m(\pi)$, Algorithm 1 implicitly requires the range of $S_H(m)/m$ is a subset of the domain of $g$, as otherwise the algorithm cannot be implemented. This constraint is naturally satisfied when $g : \mathcal{D} \to \mathbb{R}$ has domain $\mathcal{D} = \mathbb{R}^m$, such as $g(x) = e^x$, or $g(x_1, x_2) = \max\{x_1, x_2, 1\}$ and so on. However, many natural functions are not defined on the whole space, such as $g(x) = 1/x$, or $g(x_1, x_2) = x_1/x_2$. These functions arise in statistical applications such as inference on doubly-intractable problems Murray et al. [2006], Lyne et al. [2015], estimating the ratio of normalizing constants Meng and Wong [1996]. Unfortunately, Algorithm 1 cannot be implemented if $S_H(m)/m$ falls outside the domain of $g$.

Consider a concrete problem of estimating $g(m(\pi)) = 1/m(\pi)$ where $\pi$ is a probability measure on some state space $\Omega$. The problem can be naturally avoided if $S_H(m)/m \neq 0$ almost surely, which is often the case for continuous state-space $\Omega$. However, the algorithm fails for discrete state spaces. Even if $\Omega$ only contain positive numbers that are

far from 0, the resulting JOA estimator may still take 0 with positive probability, as it is constructed by the difference of two chains, see formula (3). The same problem only gets worse if the domain of $g$ is of the form $\{x \mid \|x\| \geq c\}$, where both continuous and discrete Markov chains may fail.

To address this issue, we add an extra $\delta$-transformation when necessary. Suppose $\mathcal{D} \supset \mathbb{R}^d \backslash B_\delta := \{x \mid \|x\| \geq \delta\}$ contains everything in $\mathbb{R}^d$ except for a compact set. Let $H$ be the output of the unbiased MCMC subroutine $\mathcal{S}$, then the transformation $H \to H1_{H \geq \delta} + (H + 2\delta B)1_{H < \delta}$ outputs a random variable supporting on $\mathbb{R}^d \backslash B_\delta$ while maintaining the expectation, where $B$ is uniform random variable on $\{-1, 1\}$ independent with $H$. The procedure is formally described below.

---

**Algorithm 3** The $\delta$-transformation

---

**Input:**
- A subroutine $\mathcal{S}$ for generating unbiased estimators of $m(\pi)$.
- A positive constant $\delta$.

**Output:** An unbiased estimator of $m(\pi)$ supporting on $\mathbb{R}^d \setminus B_\delta$.

Call $\mathcal{S}$ once and label the outputs by $H$.

**If** $H \geq \delta$:

    **Return** $H$.

**Else:**

    Flip a fair coin, **Return** $H + 2\delta$ if the coin comes up heads. **Return** $H - 2\delta$ otherwise.

---

The following proposition gives the validity proof of Algorithm 3.

**Proposition 2.** *Let $\tilde{H}$ be the output of Algorithm 3, then $\tilde{H} \geq \delta$ and $\mathbb{E}[\tilde{H}] = \mathbb{E}[H] = m(\pi)$.*

*Proof.* It suffices to show the unbiasedness of $\tilde{H}$. Notice that $\tilde{H} = H1_{H \geq \delta} + (H + 2\delta B)1_{H < \delta}$ where $B \sim \mathsf{U}\{-1, 1\}$ is independent with $H$. Therefore,

$$\mathbb{E}[\tilde{H}] = \mathbb{E}[H1_{H \geq \delta}] + \mathbb{E}[(H + 2\delta B)1_{H < \delta}] = \mathbb{E}[H1_{H \geq \delta}] + \mathbb{E}[H1_{H < \delta}] = \mathbb{E}[H].$$

$\square$

The $\delta$-transformation is used as a post-processing technique for the outputs of the unbiased MCMC algorithm. The whole algorithm for $g(m(\pi))$ is described in Algorithm 4 below. Algorithm 4 is the same as Algorithm 1 execpt for the post-processing procedure (Algorithm 3).

Although the $\delta$-transformation trick allows one to work with functions not defined on the entire space, the assumption $\mathcal{D} \supset \mathbb{R}^d \setminus B_\delta$ still excludes many important functions. One typical example is $g(x) = \log(x)$ which is defined on $(0, \infty)$. In order to apply Algorithm 1, it is necessary to design a non-negative unbiased estimator of $m(\pi)$, which is in general quite difficult, and sometimes impossible as discussed in Jacob and Thiery [2015]. For example, the JOA estimator cannot be directly applied even if $\Omega$ contains only positive numbers – as it is the difference between two Markov chains. The domain constraint of $g$ can be viewed as a limitation of our method, and we hope to report further progress on relaxing this constraint in future works.

---

**Algorithm 4** Unbiased Multilevel Monte-Carlo estimator after $\delta$-transformation

---

**Input:**
- A subroutine $\mathcal{S}$ for generating unbiased estimators of $m(\pi)$.
- A function $g : \mathbb{R}^d \to \mathbb{R}$.
- The parameter $p$ for geometric distribution.

**Output:** Unbiased estimator of $g(m(\pi))$.

1. Sample $N$ from the geometric distribution $\mathsf{Geo}(p)$
2. Call $\mathcal{S}$ for $2^N$ times and label the outputs by $H_1, \ldots, H_{2^N}$
3. Call Algorithm 3 for each $H_i$ and label the outputs by $\tilde{H}_1, \ldots, \tilde{H}_{2^N}$
4. Calculate the quantities $S_{\tilde{H}}(2^N)$, $S_{\tilde{H}}^{\mathsf{O}}(2^{N-1})$ and $S_{\tilde{H}}^{\mathsf{E}}(2^{N-1})$ by (4),(5)
5. Calculate $\Delta_N = g\left(S_{\tilde{H}}(2^N)/2^N\right) - \frac{1}{2}\left(g\left(S_{\tilde{H}}^{\mathsf{O}}(2^{N-1})/2^{N-1}\right) + g\left(S_{\tilde{H}}^{\mathsf{E}}(2^{N-1})/2^{N-1}\right)\right)$
   **Return:**  $W = \Delta_N/p_N + g(\tilde{H}_1)$.

---

## 3.4  Theoretical results

With all the notations as above, we are ready to state our technical assumptions and prove the theoretical results. Our theoretical analysis will focus on the unbiased estimator described in Algorithm 1. All the results still go through if the $\delta$-transformation is needed. Recall that $g$ is a function from $\mathcal{D}$ to $\mathbb{R}$, and $H_1, H_2, \cdots$ are *i.i.d.* unbiased estimators of $m(\pi)$. Now we denote by $V_n \subset \mathbb{R}^d$ the range of $(H_1 + \cdots + H_n)/n$ for every $n$ and $V := \cup_{n=1}^\infty V_n$. Our assumptions are posed on both $g$ and $H_i$:

**Assumption 3.1** (Domain)**.** *The function $g : \mathcal{D} \to \mathbb{R}$ satisfies $V \subset \mathcal{D}$. Moreover, $m(\pi)$ is in the interior of $\mathcal{D}$, i.e., $m(\pi) \in \mathcal{D}^\circ$.*

**Assumption 3.2** (Consistency)**.** $\mathbb{E}[g(\frac{S_H(n)}{n})] \to g(m(\pi))$ *as $n \to \infty$.*

**Assumption 3.3** (Smoothness)**.** *The function $g$ is continuously differentiable in a neighborhood of $m(\pi)$, and $Dg\left(\cdot\right)$ is locally Hölder continuous with exponent $\alpha > 0$. In other words, there exists $\varepsilon > 0$, $\alpha > 0$ and $c = c(\epsilon) > 0$ such that the following inequality holds for every $x, y \in (m(\pi) - \epsilon, m(\pi) + \epsilon)$:*

$$\|Dg(x) - Dg(y)\| \le c\|x - y\|^\alpha.$$

**Assumption 3.4** (Moment)**.** *There exists some $l > 2 + \alpha$ such that $H$ has finite $l$-th moments, i.e.,*

$$\mathbb{E}[\|H_1\|_l^l] = \sum_{i=1}^m \mathbb{E}[H_{1,i}^l] < \infty.$$

**Assumption 3.5** (Smoothness–Moment Tradeoff)**.** *There exist constants $s > 1$, $\alpha_s \in \mathbb{R}$, and $\mathcal{C}_s > 0$ such that $2\alpha_s + (s-1)l > 2s$ and $\mathbb{E}(|\Delta_n|^{2s}) \le \mathcal{C}_s 2^{-\alpha_s n}$ for every $n \ge 0$, where*

$$\Delta_n = \begin{cases} g\left(S_H(2^n)/2^n\right) - \frac{1}{2}\left(g\left(S_H^{\mathsf{O}}(2^{n-1})/2^{n-1}\right) + g\left(S_H^{\mathsf{E}}(2^{n-1})/2^{n-1}\right)\right) & n \ge 1 \\ g(H_1) & n = 0. \end{cases}$$

Now we briefly comment on the Assumptions 3.1 − 3.5. The descriptions below are mostly pedagogical, and the detailed proofs are deferred to the Appendix (Section A).

The Domain Assumption 3.1 guarantees Algorithm 1 can be implemented. When $g$ does not directly satisfy this assumption, but $\mathcal{D} \supset \mathbb{R}^d \setminus B_\delta$, then we apply the $\delta$-transformation (Algorithm 3 and 4) to enforce the first half of Assumption 3.1 holds. All the theoretical results still hold as long as the second half of Assumption 3.1 holds.

The consistency Assumption 3.2 is also expected and somewhat necessary. It appears in related works Vihola [2018], Blanchet and Glynn [2015], Zhou et al. [2021] explicitly or implicitly. The Law of Large Numbers guarantees $S_H(n)/n \to m(\pi)$ almost surely, therefore $g(S_H(n)/n) \to g(m(\pi))$ almost surely due to the continuity of $g$. Assumption 3.2 requires $\mathbb{E}[g(S_H(n)/n)] \to \mathbb{E}[g(m(\pi))]$, which is generally satisfied if one can use the dominated convergence theorem.

The Smoothness Assumption 3.3 guarantees both $g$ is smooth enough at a neighborhood of $m(\pi)$, and the derivative of $g$ is Hölder continuous. When $g$ is infinitely differentiable, and there is no singularity on a neighborhood of $m(\pi)$, then we expect Assumption 3.3 to hold with $\alpha \geq 1$. We emphasis that we only require $g$ to be locally Hölder continuous neary $m(\pi)$, which is much weaker than requiring $g$ to be globally Hölder continuous.

The Moment Assumption 3.4 requires more than $(2+\alpha)$-th moment of the unbiased estimator $H_i$. When the JOA estimator is used for generating $H_i$, Assumption 3.4 generally holds when $\pi$ has strictly more than $l$-th moment, and the coupling time $\tau$ has a very light tail. We recall that a $\pi$-stationary Markov chain with transition kernel $P$ is said to be geometrically ergodic if there is a $\gamma \in (0,1)$ and a function $C : \Omega \to (0, \infty)$ such that

$$\|P^n(x, \cdot) - \pi\|_{\mathsf{TV}} \leq C(x)\gamma^n,$$

for $\pi$–a.s. $x$. Geometric ergodicity is a centered notion in MCMC theory. There is a large body of literatures, including but not limited to, Mengersen and Tweedie [1996], Roberts and Tweedie [1996a,b], Jarner and Hansen [2000], Wang [2020], Livingstone et al. [2019], that shows a wide family of MCMC algorithms are geometrically ergodic.

Our result for guaranteeing Assumption 3.4 is the following.

**Proposition 3** (Verifying Assumption 3.4, informal)**.** *Suppose the Markov chain $P$ is $\pi$-stationary and geometrically ergodic, and $f$ is a measurable function with finite p-th moment under $\pi$ for any $p > l$. Suppose also there exists a set $\mathcal{S} \subset \Omega$, a constant $\tilde{\epsilon} \in (0,1)$ such that*

$$\inf_{(x,y)\in\mathcal{S}\times\mathcal{S}} \bar{P}((x,y),\mathcal{D}) \geq \tilde{\epsilon},$$

*where $\mathcal{D} := \{(x,x) : x \in \Omega\}$ is the diagonal of $\Omega \times \Omega$. Then the JOA estimator $H_k(Y,Z) := f(Y_k) + \sum_{i=k+1}^{\tau-1}(f(Y_i) - f(Z_{i-1}))$ has a finite l-th moment, and therefore satisfies Assumption 3.4.*

The formal description of the above proposition and the detailed proofs will be deferred to Appendix A.3. Proposition 3 shows the existence of existence of the $l$-th moment of the JOA estimator for $l > 2$. The proof uses very similar techniques as Jacob et al. [2020]. It can be viewed as a slightly stronger version of Proposition 3.1 in Jacob et al. [2020], where the authors established the finite second-order moment of the JOA estimator.

The Tradeoff Assumption 3.5 bounds the magnitude of $\mathbb{E}(\|\Delta_n\|^{2s})$. The condition

$$2\alpha_s + (s-1)l > 2s$$

reflects the tradeoff between the smoothness of $g$ and the moment assumption on $H_i$. Consider the following scenarios:

- Suppose $g$ is at least twice continuously differentiable, and the derivative $Dg$ is Lipschitz continuous. Then we have $\Delta_n = \mathcal{O}((S_H(2^n)/2^n)^2)$ by Taylor expansion. Meanwhile, the Central Limit Theorem (CLT) gives us $S_H(2^n)/2^n, S_H^{\mathsf{O}}(2^{n-1})/2^{n-1}$, and $S_H^{\mathsf{E}}(2^{n-1})/2^{n-1}$ are all of the magnitude $\mathcal{O}_p(2^{-n/2})$, which in turn shows $\Delta_n = \mathcal{O}_p(2^{-n})$. Therefore we expect to choose $\alpha_s = s$ and therefore Assumption 3.5 is true for positive $l$. In this case Assumption 3.5 is weaker than Assumption 3.4.
- Suppose $g$ is at most of linear growth, i.e., $|g(x)| \leq c(1 + \|x\|)$. When $g$ does not have higher-order derivatives but still have some global control on the growth rate, we can not directly use Taylor expansion to cancel out the linear terms. In this case we can only bound $\Delta_n$ by $\mathcal{O}(S_H(2^n)/2^n)$, which is $\mathcal{O}_p(2^{-n/2})$ again by the CLT. We expect to choose $\alpha_s = s/2$ and it thus requires $l > s/(s-1)$. This is also the assumption in Blanchet and Glynn [2015], Blanchet et al. [2019].
- Suppose there is no special smoothness assumption on $g$, but $\mathbb{E}[\|\Delta_n\|^{2s}]$ is uniformly bounded. Then we expect to choose $\alpha_s = 0$, and therefore $l > 2s/(s-1)$.

As we can see from the above discussions, stronger smoothness requirements on $g$ result in weaker assumptions on the moment of $H_i$, and vice versa. Theoretically, Assumption 3.5 even holds when $\alpha_s < 0$, provided that $l$ is large enough.

Our main theoretical result is as follows.

**Theorem 1.** *Under Assumption 3.1 – 3.5, let*

$$\gamma := \min\{\alpha, \frac{\alpha_s}{s} + \frac{(s-1)l}{2s} - 1\} > 0.$$

*if $N \in \{1, 2, \ldots\}$ is geometrically distributed with success parameter $p \in \left(\frac{1}{2}, 1 - \frac{1}{2^{(1+\gamma)}}\right)$, then the estimator*

$$W := \frac{\Delta_N}{p_N} + g(H_1)$$

*described in Algorithm 1 satisfies:*

1. $\mathbb{E}[W] = g(m(\pi))$,
2. *There exists a constant $C$ such that*

$$\mathsf{Var}(W) \leq \mathbb{E}[W^2] \leq Cp^{-1} \frac{2^{-(1+\gamma)}}{1 - \left((1-p)2^{1+\gamma}\right)^{-1}} < \infty.$$

3. *The expected computational cost of Algorithm 1 is finite.*

The proof of Theorem 1 relies on the following key lemma to upper bound the second moment of $\Delta_n$. The idea of bounding $\mathbb{E}[|\Delta_n|^2]$ appears in many relevant literatures such as Blanchet and Glynn [2015], Blanchet et al. [2019], Vihola [2018], Rhee and Glynn [2015] under various technical assumptions. Our tradeoff assumption 3.5 seems to be novel to our best knowledge.

**Lemma 1.** *We have*

$$\mathbb{E}[|\Delta_n|^2] = C2^{-(1+\gamma)n},$$

*where*

$$\gamma = \{\alpha, \frac{\alpha_s}{s} + \frac{(s-1)l}{2s} - 1\} > 0,$$

*and $C = C(m, l, \epsilon, s, \alpha)$ is a constant provided that Assumption 3.1 – 3.5 are satisfied.*

Though the detailed proof will be deferred to the Appendix A.2, we sketch the proof idea here to provide some insights. For now, we temporarily assume $g$ is smooth and has bounded second order derivative. Since

$$\Delta_n = g\left(S_H(2^n)/2^n\right) - \frac{1}{2}\left(g\left(S_H^O(2^{n-1})/2^{n-1}\right) + g\left(S_H^E(2^{n-1})/2^{n-1}\right)\right), \qquad (6)$$

each term in the above summation is $\mathcal{O}_p(2^{-n/2})$ by the Central Limit Theorem. Therefore, using the triangle inequality will give an $\mathcal{O}(2^{-n/2})$ upper bound for $\mathbb{E}[|\Delta_n|^2]$, which is strictly weaker than Lemma 1. The key observation is the antithetic design of $\Delta_n$ reduces the variance by a factor of $2^{-\Omega(1)n}$. More precisely, recall that

$$\frac{S_H(2^n)}{2^n} = \frac{1}{2}\left(\frac{S_H^O(2^{n-1})}{2^{n-1}} + \frac{S_H^E(2^{n-1})}{2^{n-1}}\right).$$

By Taylor expansion we have:

$$g(a) = g\left(\frac{a+b}{2}\right) + g'\left(\frac{a+b}{2}\right)\left(\frac{a-b}{2}\right) + \mathcal{O}((a-b)^2),$$

and

$$g(b) = g\left(\frac{a+b}{2}\right) + g'\left(\frac{a+b}{2}\right)\left(\frac{b-a}{2}\right) + \mathcal{O}((a-b)^2).$$

Therefore,

$$g\left(\frac{a+b}{2}\right) - \frac{1}{2}\left(g(a) + g(b)\right) = \mathcal{O}((a-b)^2),$$

so the anthithesic difference cancels out the constant and linear terms in the Taylor expansion, leaving the second order term as the dominating term. Taking $a = \frac{S_H^O(2^{n-1})}{2^{n-1}}$ and $b = \frac{S_H^E(2^{n-1})}{2^{n-1}}$, we know

$$\Delta_n = \mathcal{O}\left(\left(\frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}}\right)^2\right) = \mathcal{O}_p(2^{-n}).$$

Therefore

$$\mathbb{E}[\Delta_n^2] = \mathcal{O}(2^{-2n}),$$

corresponding to the case $\gamma = 1$ in Lemma 1.

This gives the intuition of the proof idea. Our real technical assumptions (3.3) are more general than the idealized assumption above. It only requires $\alpha$-Hölder continuity of $Dg$ on a neighborhood of $m(\pi)$. Therefore, we discuss the behavior of $|\Delta_n|^2$ when $\frac{S_H(2^n)}{2^n}$ is closed to, or far away from its expectation. In both cases, we show the expected value of $|\Delta_n|^2$ is $\mathcal{O}(2^{-(1+\Omega(1))n})$. The details of the proof can be found in Appendix A.2.

With Lemma 1 in hand, we are ready to show Theorem 1.

*Proof of Theorem* 1. We will first show Statement 1 assuming Statement 2 holds. Then we show both Statement 2 and 3 holds. The proof of Statement 2 depends heavily on Lemma 1.

*Proof of Statement* 1: Suppose $W$ has a finite second moment, then the conditional expectation $\mathbb{E}[W \mid N]$ is well defined (see Section 4.1 of Durrett [2019] for details). The law of iterated expectation yields

$$\mathbb{E}[W] = \mathbb{E}\big[\mathbb{E}[W \mid N]\big] = \mathbb{E}[g(H_1)] + \mathbb{E}\Big[\frac{\mathbb{E}[\Delta_n \mid N]}{p_N}\Big] = \mathbb{E}[g(H_1)] + \mathbb{E}\Big[\frac{d_N}{p_N}\Big],$$

where $d_n = \mathbb{E}[g(S_H(2^n)/2^n)] - \mathbb{E}[g(S_H(2^{n-1})/2^{n-1})]$. We can further calculate $\mathbb{E}\big[\frac{d_N}{p_N}\big]$:

$$\mathbb{E}\Big[\frac{d_N}{p_N}\Big] = \sum_{i=1}^{\infty} \frac{d_i}{p_i} p_i = \sum_{i=1}^{\infty} d_i.$$

Therefore

$$\mathbb{E}[W] = \lim_{n\to\infty} \mathbb{E}[g(S_H(2^n)/2^n)] = g(m(\pi)),$$

as desired. The last equality uses Assumption 3.2.

*Proof of Statement* 2: Now we show $\mathbb{E}[W^2] < \infty$. We can directly calculate:

$$\mathbb{E}[W^2] \leq 2\Big(\mathbb{E}[g(H_1)^2] + \mathbb{E}\Big[\frac{\Delta_N^2}{p_N^2}\Big]\Big).$$

It suffices to show $\mathbb{E}\big[\frac{\Delta_N^2}{p_N^2}\big] < \infty$. We have

$$\mathbb{E}\Big[\frac{\Delta_N^2}{p_N^2}\Big] = \sum_{n=1}^{\infty} \frac{\mathbb{E}[\Delta_n^2]}{p_n} = \sum_{n=1}^{\infty} \mathbb{E}[\Delta_n^2](1-p)^{-n+1}p^{-1}.$$

By Lemma 1,

$$\mathbb{E}\Big[\frac{\Delta_N^2}{p_N^2}\Big] \leq Cp^{-1}(1-p)\sum_{n=1}^{\infty} 2^{-(1+\gamma)n}(1-p)^{-n} \tag{7}$$

$$= Cp^{-1}(1-p)\sum_{n=1}^{\infty} \Big((1-p)2^{1+\gamma}\Big)^{-n} \tag{8}$$

$$= Cp^{-1}\frac{2^{-(1+\gamma)}}{1 - \big((1-p)2^{1+\gamma}\big)^{-1}} < \infty, \tag{9}$$

where the last inequality follows from $(1-p) > 2^{-(\gamma+1)}$.

*Proof of Statement* 3: We are now in the position of bounding the computation cost of Algorithm 1. Let $C_H$ be the computation cost for implementing the unbiased MCMC subroutine $\mathcal{S}$ once. It is shown in Jacob et al. [2020] that $C_H < \infty$. The computation cost for implementing Algorithm 1 essentially comes from $2^N$ calls of the subroutine $\mathcal{S}$, where $N \sim \mathsf{Geo}(p)$. Therefore it suffices to show $2^N$ has a finite expectation. We calculate

$$\mathbb{E}[2^N] = \sum_{n=1}^{\infty} 2^n p(n) = \sum_{n=1}^{\infty} 2^n (1-p)^{n-1}p = \frac{2p}{2p-1} < \infty, \tag{10}$$

where the last inequality follows from $p > \frac{1}{2}$. $\qquad\square$

The proof of Theorem 1 also suggests the following strategy to choose the success parameter $p$. We follow the definition in Glynn and Whitt [1992], Blanchet and Glynn [2015] and define the work-normalized variance $\tilde{\sigma}^2(W)$ to be the product of the computation cost and the variance of an individual estimator. Then it follows from (7) and (10) that the work-normalized variance of $W$ is upper bounded by a constant multiple of:

$$\sum_{n=1}^{\infty} \left( (1-p)2^{1+\gamma} \right)^{-n} \times \sum_{n=1}^{\infty} \left( 2(1-p) \right)^n. \tag{11}$$

By Cauchy-Schwarz inequality, formula (11) can be minimized by choosing $p = 1 - 2^{-1-\frac{\gamma}{2}}$. When $\gamma = 1$, the parameter $p$ can be chosen as $p = 1 - 2^{-\frac{3}{2}} \approx 0.646$, recovering the result in Blanchet and Glynn [2015]. We will see more discussions on empirical strategies of choosing $p$ in Section 4.

Finally, we present two Central Limit Theorems (CLTs) of our estimator. These results follow directly from the standard Central Limit Theorem arguments from Glynn and Heidelberger [1991], Blanchet and Glynn [2015]. These results show our estimator has the optimal 'square-root' convergence rate. Confidence intervals can also be established using these CLTs.

- When the number of estimators $W_1, W_2, \ldots, W_n, \ldots$ in Algorithm 1 goes to infinity, we have
$$\left( \frac{\sum_{i=1}^n W_i}{\sqrt{n}} - g(m(\pi)) \right) \to \mathsf{N}(0, \mathsf{Var}(W_1)) \qquad \text{as } n \to \infty.$$

- Given a fixed computational budget $b$, let $N(b)$ be the number of *i.i.d.* estimators $W_1, W_2, \ldots, W_{N(b)}$ that can be generated by time $b$. Then we have
$$\sqrt{b} \cdot \left( \frac{\sum_{i=1}^{N(b)} W_i}{N(b)} - g(m(\pi)) \right) \to \mathsf{N}(0, \tilde{\sigma}^2(W)) \qquad \text{as } b \to \infty.$$

# 4 Numerical examples

In this section, we investigate the empirical performance of the proposed method with several examples. We first test the algorithm on a multivariate distribution on the high-dimensional cube to show the correctness of our estimator and study the algorithm's sensitivity to the parameter $p$ in Algorithm 1. Then we implement our algorithm on a 2-D Ising model with periodic boundaries, and show numerical results of two statistics of interest.

## 4.1 Product of inverse expectations

We begin with a toy model to illustrate the performance of our method. Let the multivariate random variable be $X = (X_1, \cdots, X_K) \in [0,1]^K$ with independent components $X_i \sim \mathsf{Beta}(i, 1)$. We are interested in estimating the product of the inverse expectation:

$$g_K \left( \mathbb{E}[X_1], \cdots, \mathbb{E}[X_K] \right) = \prod_{i=1}^{K} \frac{1}{\mathbb{E}[X_i]}$$

for $K = 1, \cdots, d$. Standard calculation yields

$$g_K \left( \mathbb{E}[X_1], \cdots, \mathbb{E}[X_K] \right) = \prod_{i=1}^{K} \frac{i+1}{i} = K + 1.$$

In other words, the ground-truth of $g_K$ grows linearly with $K$. On the other hand, $g_K$ cannot be directly expressed as an expectation of the underlying probability measure, and therefore the standard methods fail to provide unbiased estimators of this quantity.

We apply our method to this problem. We first test the sensitivity of Algorithm 1 to the parameter $p$, the success probability of the geometric distribution. Setting $K = 8$, and using the R package 'unbiasedmcmc' in Jacob et al. [2020] for estimating $\mathbb{E}[X_i]$ [2], we generate $10^5$ unbiased estimates of $g_K(\cdot)$ using Algorithm 1 with parameter $p$ ranging from 0.6 to 0.8. Figure 2 reports the empirical average and variance of the $10^5$ estimators for each $p$. As shown Figure 2a, the estimates are all accurate up to two decimal places and vary little for different $p$. To ensure both small estimating error and efficient computation, we set $p = 0.7$ in this example.
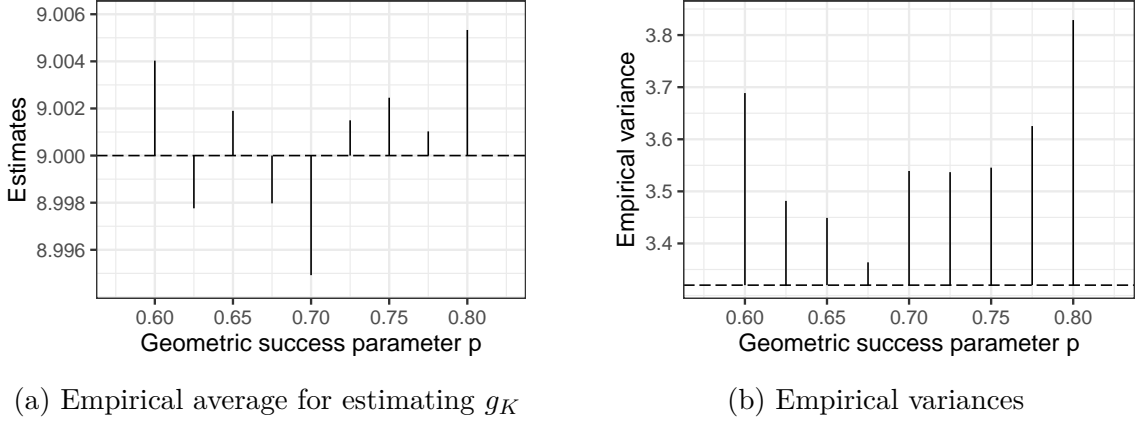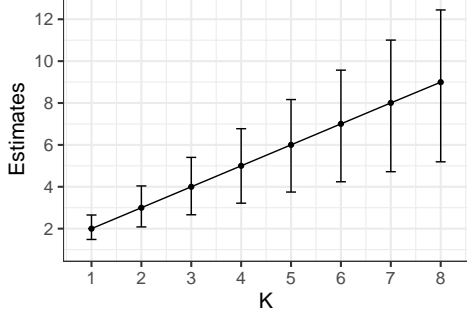


| (a) Empirical average for estimating $g_K$ | (b) Empirical variances |

Figure 2: Results for $K = 8$ with different parameters. Left: The histogram of empirical averages of $10^5$ estimators with different $p$. The solid horizontal line stands for the ground truth. Left: The histogram of empirical variances of $10^5$ estimators with different $p$.

After setting the parameter $p$, our results for $K \in \{1, 2, \cdots, 8\}$ are presented in Figure 3. For each $K$, we implement Algorithm 1 for $10^5$ times independently to generate unbiased estimators of $g_K$. The 95% Confidence Intervals (CIs) are constructed by the 2.5% and 97.5% quantiles of the $10^5$ estimators. Our point estimates and the corresponding CIs are reported in Figure 3a. As shown clearly in the figure, the point estimates are highly accurate and fit the ground-truth $g_K = K + 1$ almost perfectly. The confidence intervals get wider when $K$ increases, indicating a higher uncertainty under higher dimensionality.

To better understand the scalability, we plot the empirical variances of our estimator with the dimensionality $K$ in Figure 3b. As expected, the empirical variance grows as $K$ grows, but the correct scaling between variance and dimensionality is still unclear from our plot, leading to an interesting problem that may be worth investigating theoretically.

---

[2]In this specific example, the Beta distribution can be perfectly sampled, and there is no need to use the JOA estimator in practice. However, for illustrating our general framework, we still implement the JOA estimators for estimating $\mathbb{E}[X_i]$ via couplings of MCMC algorithms.

(a) Estimates and the 95% Confidence Intervals      (b) Empirical variances

Figure 3: Estimates of $g_K$ for $K \in \{1, 2, \cdots, 8\}$ with $p = 0.7$. Left: Estimates and the corresponding 95% Confidence Intervals from $10^5$ estimators generated by Algorithm 1. Left: The empirical variances of our estimates.

## 4.2    Ising model

We examine our method on the 2-d square-lattice Ising model. Let $\Lambda$ be a set of $n \times n$ lattice sites with periodic boundary conditions. A spin configuration $\sigma \in \{-1, 1\}^{n \times n}$ is an assignment of spins to all the lattice vertices. A 2-d Ising model is a probability distribution over all the spin configurations, defined as:

$$p_\theta(\sigma) = \frac{\exp(-\theta H(\sigma))}{Z(\theta)},$$

Here $H(\sigma) = -\sum_{\langle I,J \rangle} \sigma_i \sigma_j$ is referred to as the 'the Hamiltonian function', the sum is over all pairs of neighboring lattice sites. The normalizing constant $Z(\theta) = \sum_\sigma \exp(-\theta H(\sigma))$ is referred to as the partition function. The parameter $\theta \geq 0$ is often interpreted as the inverse temperature in statistical physics.

### 4.2.1    Inverse of natural statistics

Let us denote the 'natural statistics' of the Ising model by $h(\sigma) := -H(\sigma)$. In this example we are interested in estimating $1/\mathbb{E}_\theta[h(\sigma)]$. Standard calculation in exponential families yields:

$$\frac{1}{\mathbb{E}_\theta[h(\sigma)]} = \frac{1}{\log(Z(\theta))'} = \frac{Z(\theta)}{Z'(\theta)}.$$

Following the setups in Jacob et al. [2020], we set $n = 32$ (which means the sample space is of dimension $32^2 = 1024$) and use the JOA estimator for unbiased estimation of $\mathbb{E}_\theta[h(\sigma)]$ by coupling two single-site Gibbs samplers, and feed these estimators as inputs for the unbiased MLMC estimator with parameter $p = 0.7$, as described in Algorithm 4. We implement our estimator for a grid of $\theta$ values ranging from 0.23 to 0.40. For each $\theta$, we generate $10^5$ unbiased estimators and report our results in Figure 4 below. Similar to the observations in Jacob et al. [2020], the meeting time increases exponentially as $\theta$ increases. Therefore it may be computationally demanding to generate unbiased estimators when $\theta$ is close to its critical temperature. Meanwhile, the standard deviation has an interesting

$U$-shape pattern as $\theta$ increases, as shown in Figure 4b. We have no idea how to explain this phenomenon theoretically.



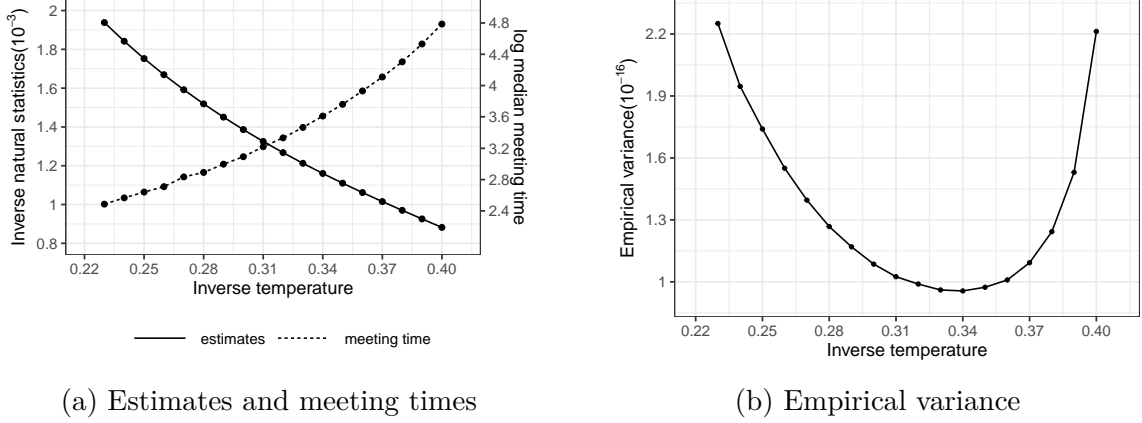(a) Estimates and meeting times        (b) Empirical variance

Figure 4: Estimates, meeting times and standard deviations of $1/\mathbb{E}_\theta[h(\sigma)]$ for $\theta \in \{0.23, 0.24, \ldots, 0.4\}$. Left: the solid line stands for the empirical averages of $10^5$ unbiased estimators from Algorithm 4. The dashed line stands for the log median meeting time of the JOA estimators.

### 4.2.2 Ratio of normalizing constants

Now we consider a more challenging task – estimating the ratio of normalizing constant $Z(\theta_1)/Z(\theta_2)$. The problem, also known as estimating the free energy differences, is of great interest in the computational physics and the statistics community. We refer the readers to Bennett [1976], Diaconis and Wang [2018], and Meng and Wong [1996] for further discussions. Since the Ising model is generally computationally intensive to be sampled perfectly (see Propp and Wilson [1996]), unbiased estimators of $Z(\theta_1)/Z(\theta_2)$ is generally unavailable for moderate $n$ in the previous literature.

We will use our method to construct unbiased estimators of $Z(\theta_1)/Z(\theta_2)$. First, we observe the following identity:

$$\sum_\sigma p_\theta(\sigma) \cdot \frac{1}{\exp(\theta h(\sigma))} = \sum_\sigma \frac{\exp(\theta h(\sigma))}{Z(\theta)} \cdot \frac{1}{\exp(\theta h(\sigma))} = \sum_\sigma \frac{1}{Z(\theta)} = \frac{2^{n^2}}{Z(\theta)}.$$

Therefore $Z(\theta) = 2^{n^2}/\mathbb{E}_\theta[e^{\theta H(\sigma)}]$. Thus the ratio can be written as

$$\frac{Z(\theta_1)}{Z(\theta_2)} = \frac{\mathbb{E}_{\theta_2}[e^{\theta_2 H(\sigma)}]}{\mathbb{E}_{\theta_1}[e^{\theta_1 H(\sigma)}]}. \tag{12}$$

Therefore, for fixed $\theta_1, \theta_2$, we call the JOA estimators for unbiased estimation for both the numerator and the denominator of (12), and feed them into Algorithm 4 for unbiased estimators of $Z(\theta_1)/Z(\theta_2)$. The JOA estimators by coupling two Gibbs samplers, using the package 'unbiasedmcmc' in Jacob et al. [2020].

We implement our method using the parameters $n = 12, p = 0.7, \theta_1 \in \{0.05, 0.07, \ldots, 0.20\}$ and $\theta_2 \in \{0.05, 0.10, 0.15, 0.20\}$ on a 500-core CPU-based computer cluster. For each

combination of $(\theta_1, \theta_2)$, we generate 200 unbiased estimators on each cluster, resulting in $10^5$ unbiased estimators in total. We present results of the normalizing constant ratio in Figure 5. The solid lines represent our estimates for $Z(\theta_1)/Z(\theta_2)$. To check the accuracy of our method, we run the Gibbs sampler 500 times for estimating both $\mathbb{E}_{\theta_1}[e^{\theta_1 H(\sigma)}]$ and $\mathbb{E}_{\theta_2}[e^{\theta_2 H(\sigma)}]$, then plug in both estimates into (12) to estimate $Z(\theta_1)/Z(\theta_2)$. It is worth mentioning that the Gibbs sampler is easy to implement, but generally biased. Our results using the plug-in estimators from the Gibbs sampler are presented in the dashed line of Figure 5. As clearly shown in the figures, both lines are extremely close to each other, indicating that our point estimates match well with the standard estimators, while retaining the unbiasedness.



(a) $Z(\theta_1)/Z(0.05)$

(b) $Z(\theta_1)/Z(0.10)$

(c) $Z(\theta_1)/Z(0.15)$
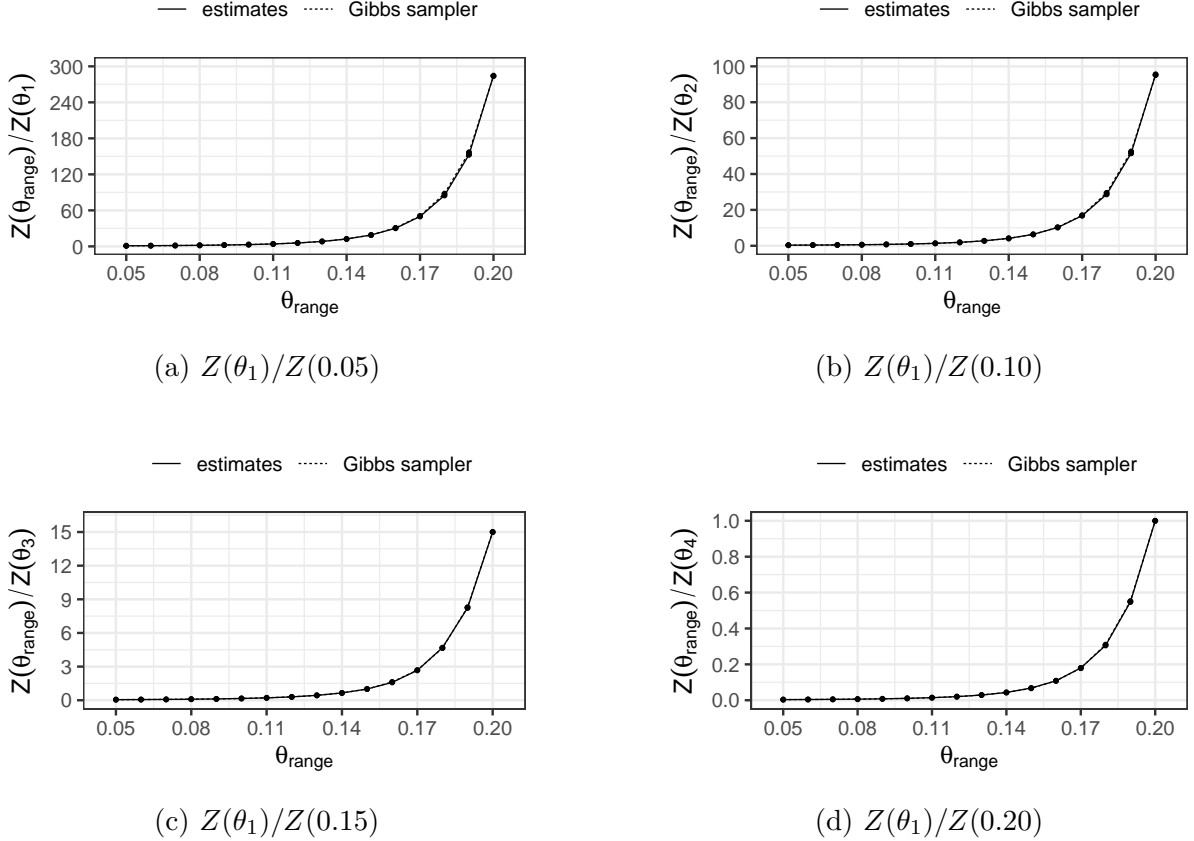
(d) $Z(\theta_1)/Z(0.20)$

Figure 5: Ratio of normalizing constant $Z(\theta_1)/Z(\theta_2)$ for $n = 12$. Solid lines represent our estimators and dash lines are estimates generated by the Gibbs sampler.

We further compare the 95% confidence intervals constructed by our Monte Carlo estimator and the classical Gibbs sampler. The Monte Carlo confidence interval is constructed by the 2.5% and 97.5% quantiles of the $10^5$ estimators, which can be justified by the CLTs in Section 3. The Gibbs sampler confidence interval is constructed using the Delta method, see Appendix B for details. Our results are reported in Table 1 and Table 2 respectively. As shown in the tables, the Monte Carlo CIs are comparable, though wider than the Gibbs sampler CIs. Therefore, our method may be less efficient than the Gibbs sampling in terms of a single estimator, which can be viewed as the price of being un-biasedness. However, there are two advantages of the Monte Carlo CIs. Firstly, thanks to the potential of parallel implementation, our method can be easily computed on large computer clusters, where the consistency occurs in the number of replications rather than

the number of iterations. In contrast, the Gibbs sampler CIs are justified by the Markov chain Central Limit Theorem [Jones, 2004], which is known to be challenging to verify [Rosenthal, 2017, Jiang et al., 2020], and is only valid when the sequential sampling algorithm is run long enough. Secondly, our Monte Carlo CIs does not require any knowledge on the gradient of $g$, while the Gibbs sampler CIs use the delta method and thus require the gradient information.

| | MLMC-MCMC | | | |
| | Z(0.05) | Z(0.10) | Z(0.15) | Z(0.20) |
| --- | --- | --- | --- | --- |
| Z(0.05) | | (0.283, 0.407) | (0.032, 0.100) | (0.0010, 0.0118) |
| Z(0.06) | (1.086, 1.262) | (0.331, 0.477) | (0.038, 0.117) | (0.0012, 0.0138) |
| Z(0.07) | (1.288, 1.544) | (0.396, 0.580) | (0.045, 0.142) | (0.0014, 0.0168) |
| Z(0.08) | (1.560, 1.955) | (0.483, 0.726) | (0.055, 0.177) | (0.0017, 0.0211) |
| Z(0.09) | (1.925, 2.563) | (0.603, 0.943) | (0.069, 0.228) | (0.0021, 0.0271) |
| Z(0.10) | (2.399, 3.483) | | (0.088, 0.305) | (0.0026, 0.0364) |
| Z(0.11) | (2.988, 4.918) | (0.959, 1.770) | (0.113, 0.419) | (0.0031, 0.0502) |
| Z(0.12) | (3.701, 7.215) | (1.187, 2.570) | (0.140, 0.601) | (0.0032, 0.0728) |
| Z(0.13) | (4.343, 11.021) | (1.406, 3.890) | (0.164, 0.886) | (0.0019, 0.1096) |
| Z(0.14) | (4.559, 17.628) | (1.441, 6.183) | (0.161, 1.374) | (-0.0037, 0.1707) |
| Z(0.15) | (3.256, 29.432) | (0.986, 10.230) | | (-0.0182, 0.2718) |
| Z(0.16) | (-1.795, 51.539) | (-0.877, 17.926) | (-0.274, 3.721) | (-0.0589, 0.4605) |
| Z(0.17) | (-18.290, 95.085) | (-6.264, 32.607) | (-1.259, 6.617) | (-0.1466, 0.8022) |
| Z(0.18) | (-61.301, 183.525) | (-20.757, 62.913) | (-3.846, 12.488) | (-0.4067, 1.4639) |
| Z(0.19) | (-183.972, 376.751) | (-61.099, 128.451) | (-11.053, 24.791) | (-1.0822, 2.8941) |
| Z(0.20) | (-502.632, 814.170) | (-180.765, 279.066) | (-30.252, 52.056) | |

Table 1: 95% Confidence Interval of MLMC-MCMC estimator

| | **Gibbs Sampling** | | | |
| | Z(0.05) | Z(0.10) | Z(0.15) | Z(0.20) |
|---|---|---|---|---|
| Z(0.05) | | (0.316, 0.332) | (0.047, 0.063) | (0.0018, 0.0049) |
| Z(0.06) | (1.173, 1.197) | (0.375, 0.393) | (0.056, 0.075) | (0.0021, 0.0058) |
| Z(0.07) | (1.402, 1.437) | (0.448, 0.472) | (0.067, 0.090) | (0.0025, 0.0069) |
| Z(0.08) | (1.750, 1.806) | (0.561, 0.592) | (0.084, 0.112) | (0.0031, 0.0087) |
| Z(0.09) | (2.204, 2.289) | (0.707, 0.749) | (0.106, 0.142) | (0.0039, 0.0110) |
| Z(0.10) | (3.013, 3.159) | | (0.146, 0.195) | (0.0054, 0.0151) |
| Z(0.11) | (3.821, 4.157) | (1.231, 1.354) | (0.187, 0.253) | (0.0070, 0.0195) |
| Z(0.12) | (5.421, 5.973) | (1.748, 1.944) | (0.267, 0.362) | (0.0099, 0.0278) |
| Z(0.13) | (7.520, 8.384) | (2.426, 2.728) | (0.372, 0.506) | (0.0138, 0.0389) |
| Z(0.14) | (11.353, 13.303) | (3.667, 4.323) | (0.570, 0.791) | (0.0213, 0.0604) |
| Z(0.15) | (15.530, 20.702) | (5.023, 6.718) | | (0.0304, 0.0896) |
| Z(0.16) | (21.917, 31.041) | (7.090, 10.071) | (1.135, 1.788) | (0.0437, 0.1318) |
| Z(0.17) | (39.114, 59.134) | (12.656, 19.181) | (2.037, 3.386) | (0.0791, 0.2464) |
| Z(0.18) | (89.273, 109.956) | (28.857, 35.703) | (4.529, 6.468) | (0.1707, 0.4895) |
| Z(0.19) | (116.169, 195.468) | (37.598, 63.388) | (6.093, 11.109) | (0.2396, 0.7929) |
| Z(0.20) | (159.443, 444.158) | (51.615, 143.982) | (8.451, 24.868) | |

Table 2: 95% Confidence Interval of Gibbs Sampling estimator

## 4.3 Nested expectation

Finally, we implement our estimator for estimating the following nested expectation

$$U := \mathbb{E}_{\theta_1}[\max_d \mathbb{E}_{\theta_2|\theta_1}[f_d(\theta_1, \theta_2)|\theta_1]],$$

where the set of functions $\{f_d\}_{d=1}^D$ and the joint distribution over $\theta_1, \theta_2$ will be clear soon. The quantity $\max_d \mathbb{E}_{\theta_2|\theta_1}[f_d(\theta_1, \theta_2)|\theta_1$ is often interpreted as the utility or the optimal outcome over $D$ possible choices given the information of $\theta_1$. Estimating $U$ is often of interest in sequential decision making, or modular inference problems. Since $U$ contains a nested expectation, with an out expectation over $\theta_1$ and an inner expectation over $\theta_2|\theta_1$, the vanilla Monte Carlo approach (sample $N_1$ realizations of $\theta_1$, and sample $N_2$ realizations of $\theta_2$ given each $\theta_1^{(i)}$) typically has suboptimal computational complexity $\mathcal{O}(\epsilon^{-3})$ or even $\mathcal{O}(\epsilon^{-4})$ for $\epsilon$ root mean square error (rMSE) under varying assumptions Giles and Goda [2019]. Therefore, MLMC methods have been proposed when both $\theta_1$ and $\theta_2|\theta_1$ can be perfectly sampled. The case where $\theta_2|\theta_1$ can only be approximately sampled is still considered as an open problem.

We construct an unbiased estimator of $U$ using the methodology described in Section 3.2. In this example, suppose we have two models. The first model comprises parameter $\theta_1$ with prior $\pi_1(\theta_1)$, data $Y_1$ with likelihood $p_1(y|\theta_1)$, the second model comprises parameter $\theta_2$ with prior $\pi_2(\theta_2)$, data $Y_2$ with likelihood $p_2(y|\theta_1, \theta_2)$. The cut distribution is defined as:

$$\pi^\star(\theta_1, \theta_2) := \pi(\theta_1|Y_1)\pi(\theta_2|Y_2, \theta_1). \tag{13}$$

This is different from the usual posterior distribution

$$\pi(\theta_1, \theta_2 | Y_1, Y_2) = \pi(\theta_1 | Y_1, Y_2)\pi(\theta_2 | Y_2, \theta_1). \tag{14}$$

In the cut model, the distribution of $\theta_1$ depends on the observations from the first model ($Y_1$) but not the observations from the second model ($Y_2$). In contrast, the posterior distribution of $\theta_1$ depends on both $Y_1$ and $Y_2$. Since the cut model prevents the information in the second model from influencing the inference on the first, it is often used as an alternative to Bayes full posterior in the presence of model misspecification. The investigation of the cut model is an active research area, we refer the readers to [Plummer, 2015, Jacob et al., 2020, Pompe and Jacob, 2021] for more discussions.

Conducing inference on the cut model is challenging. The conditional distribution $\pi(\theta_2 | Y_2, \theta_1)$ is usually only known up a normalizing constant $Z(\theta_1)$ which depends on $\theta_1$. Standard MCMC methods on the joint space $(\theta_1, \theta_2)$ cannot be directly implemented due to the intractability of $Z(\theta_1)$.

Going back to our case, we consider the real-data example used in [Plummer, 2015, Jacob et al., 2020] from epidemiology, which is motivated from a study of the international correlation between human papilloma virus (HPV) prevalence and cervical cancer incidence Maucort-Boulch et al. [2008]. The first module consists of high-risk HPV prevalence data from 13 countries. The data $Y_1 = \{(Z_i, N_i)\}_{i=1}^{13}$ consists 13 pair of integers, where $Z_i$ is the number of women infected with HPV, and $N_i$ is the with population size of country $i$. We assume a beta prior $\mathsf{Beta}(1,1)$ on each component of $\theta_1 = (\theta_{1,1}, \ldots, \theta_{1,13})$ independently, and an independent binomial likelihood $Z_i \sim \mathsf{Binom}(N_i, \theta_i)$ for each $i$. This yields a product beta posterior for $\theta_1$.

The second module consists of the cancer data from the same 13 countries. The data $Y_2 = \{(X_{1,i}, X_{2,i})\}_{i=1}^{13}$ consists 13 pair of integers, where $X_{1,i}$ is numbers of cancer cases arising from $X_{2,i}$ woman-years of follow-up. We assume a bivariate normal prior with mean vector 0 and a diagonal covariance matrix with variance $10^3$ per component on the parameter $\theta_2 = (\theta_{2,1}, \theta_{2,2}) \in \mathbb{R}^2$, and a Poisson regression model:

$$X_{1,i} \sim \mathsf{Poi}(\exp(\lambda_i)),$$

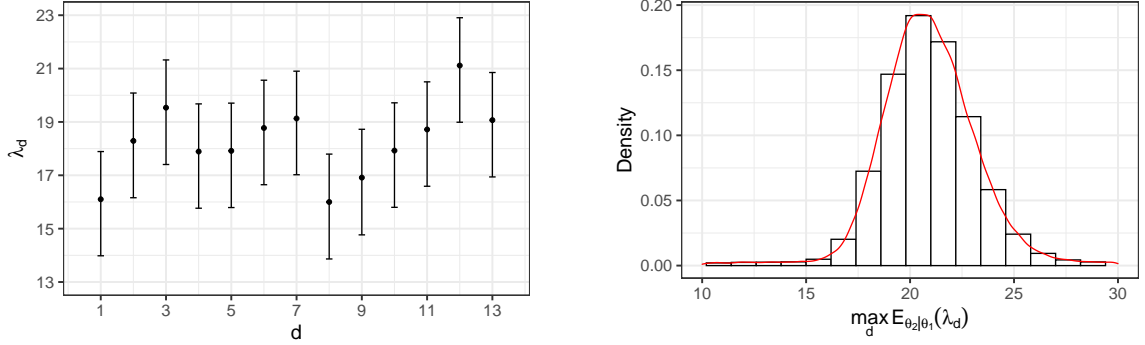where $\lambda_i = \theta_{2,1} + \theta_{1,i}\theta_{2,2} + X_{2,i}$.

Under the cut model, the first parameter $\pi(\theta_1 | Y_1)$ can be directly sampled from the product beta distribution, and the second parameter can be approximately sampled from $\pi(\theta_2 | Y_2, \theta_1)$ using MCMC. Suppose we are interested in estimating:

$$U := \mathbb{E}_{\theta_1}[\max_{d \in \{1,2,\ldots,13\}} \mathbb{E}_{\theta_2 | \theta_1}[\lambda_d]],$$

which corresponds to the expectation of largest parameter in the Poisson regression after observing $\theta_1 \sim \pi(\theta_1 | Y_1)$. Notice that though the JOA estimator can be used to estimate each $\mathbb{E}[\lambda_d]$ unbiasedly under $\pi(\theta_2 | Y_2 \theta_1)$ for every fixed $\theta_1$, directly taking the maximum over the JOA estimator yields a biased estimator of $U$ due to

$$\mathbb{E}[\max\{X_1, \ldots, X_d\}] \geq \max\{\mathbb{E}[X_1], \ldots, \mathbb{E}[X_d]\}.$$

We implement Algorithm 2 with parameter ($p = 0.7$) to get unbiased estimators of $U$. In each run, we first sample one $\theta_1$ from the product beta posterior, and then use the JOA estimator by using the R package 'unbiasedMCMC' Jacob et al. [2020] to generate estimators of $\mathbb{E}_{\theta_2|\theta_1}[\lambda_d]$ for each $d$. Finally we use the unbiased MLMC method on the JOA estimators to eliminate the bias. Our estimates are presented in Figure 6 below. The left subplot gives the estimates and their CIs of $\lambda_d$ for each $d$. The right subplot gives the histogram and the fitted curve from $10^5$ unbiased estimators of $U$. The left subplot suggests the 12-th country in our data set has the largest parameter $\lambda_d$, which is around 21. This is consistent with the result from our unbiased estimator on the right subplot.



(a) Estimates and confidence intervals for $\lambda_d$, computed from $10^5$ JOA estimators.  (b) Histogram of $\mathbb{E}_{\theta_1}[\max_{d \in \{1,2,\ldots,13\}} \mathbb{E}_{\theta_2|\theta_1}[\lambda_d]]$ computed from $10^5$ calls of Algorithm 2.

Figure 6

# 5  Future works

Based on the combination of the unbiased MCMC and MLMC method, we propose unbiased estimators of $g(\mathbb{E}_\pi[f])$ when $\pi$ can only be approximately sampled from MCMC methods. Although promising, the existing framework (Algorithm 1 and its variants) is still not flexible enough to handle many important applications in practice. Let us get back to the three examples introduced at the beginning of Section 1. We highlight the following directions that we hope to investigate further:

- *Unbiased estimation for general $\mathcal{T}$*: In the existing works, $\mathcal{T}$ is assumed to be a function of the expectation, i.e., $\mathcal{T}(\pi) := g(\mathbb{E}_\pi[f])$. This precludes many important applications where $\mathcal{T}$ depends directly on the probability measure, instead of the expectation of some probability measure. Examples include the quantile estimation ($\mathcal{T}$ is a quantile of a measure), maximum a posteriori (MAP) estimation ($\mathcal{T}$ is the maximum of a probability density/mass function), and many optimization problems ($\mathcal{T}$ is the minimum of the loss function). We plan to either develop a more general method to include some/all the applications above, or investigate the possibility of converting these problems into easier ones under suitable assumptions.
- *The domain problem*: Taking a step back, many challenges remain even assuming $\mathcal{T}(\pi) := g(\mathbb{E}_\pi[f])$. Algorithm 1 implicitly requires the range of $S_H(m)/m$ is a subset of the domain of $g$, as otherwise the algorithm cannot be implemented. For

example, our algorithm fails when $g(x) = \sqrt{x}$ which is defined on the positive half-line. This is because the JOA estimator may not always be non-negative, though it has a non-negative expectation. We want to tackle the domain problem in this project. As remarked by several authors [Lyne et al., 2015, Jacob and Thiery, 2015], the domain problem is deeply connected with the sign problem in computational physics, which is NP-hard in its general form. Research progress on the domain problem should not only let us improve our existing framework, but also benefit both the statistics and physics communities.

- *Improved unbiased estimator design*: The efficiency of the existing estimator (Algorithm 1) is still pretty much unexplored. Our theoretical result (Theorem 1) shows the unbiasedness, finite variance, and finite computational cost of the proposed estimator. In practice, however, we find the implementation time can be slow when the dimension is high or when the Markov chain mixes slowly. In particular, empirical results suggest the parameter $p$ for the geometric distribution in Algorithm 1 has a significant influence on both the variance and the computation cost. Therefore, it is an interesting problem to find the optimal parameter and the tradeoff between the computational and statistical efficiency.

# A  Proofs

## A.1  Auxiliary Lemmas

In this section we prove some auxiliary results that will be used throughout the technical proofs. We start (without proof) the well-known Marcinkiewicz-Zygmund inequality, and then prove two useful corollaries based on this inequality.

**Lemma 2** (Marcinkiewicz-Zygmund inequality [Marcinkiewicz and Zygmund, 1937]). *If $X_1, \cdots, X_n$ are independent random variables with $\mathbb{E}[X_i] = 0$ and $\mathbb{E}\left[|X_i|^p\right] < \infty$ for some $p > 2$. Then,*

$$\mathbb{E}\left[\left|\sum_{i=1}^n X_i\right|^p\right] \leq C_p \mathbb{E}\left[\left(\sum_{i=1}^n |X_i|^2\right)^{p/2}\right],$$

*where $C_p$ is a constant that only depends on $p$.*

One corollary of the Marcinkiewicz-Zygmund inequality is:

**Corollary 1.** *With all the assumptions as above, if we further assume that $X_1, \cdots, X_n$ are i.i.d. . Then,*

$$\mathbb{E}\left[\left|\frac{1}{n}\sum_{i=1}^n X_i\right|^p\right] \leq C_p \frac{\mathbb{E}|X_1|^p}{n^{p/2}}$$

*for every $p \geq 2$.*

*Proof of Corollary 1.* Applying the Marcinkiewicz-Zygmund inequality on $(X_1/n, X_2/n, \ldots, X_n/n)$, we have:

$$\mathbb{E}\left[\left|\frac{1}{n}\sum_{i=1}^n X_i\right|^p\right] \leq C_p \mathbb{E}\left[\left(\sum_{i=1}^n \left|\frac{X_i}{n}\right|^2\right)^{p/2}\right].$$

Since $x^{p/2}$ is convex, we have

$$\left(\sum_{i=1}^n \left|\frac{X_i}{n}\right|^2\right)^{p/2} = \left(\frac{1}{n}\sum_{i=1}^n \frac{|X_i|^2}{n}\right)^{p/2} \leq \frac{1}{n}\sum_{i=1}^n \frac{|X_i|^p}{n^{p/2}}.$$

Taking expectation on both sides of the above inequality yields

$$\mathbb{E}\left[\left(\sum_{i=1}^n \left|\frac{X_i}{n}\right|^2\right)^{p/2}\right] \leq \frac{\mathbb{E}|X_1|^p}{n^{p/2}},$$

and our desired inequality follows. $\qquad\square$

The Marcinkiewicz-Zygmund inequality naturally generalizes to random vectors.

**Corollary 2** (Multivariate Marcinkiewicz-Zygmund inequality). *Let $X_1, \cdots, X_n$ be i.i.d. random vectors in $\mathbb{R}^m$, with $\mathbb{E}[X_i] = \mathbf{0}$ and $\mathbb{E}[\|X_i\|_p^p] = \mathbb{E}[\sum_{j=1}^m |X_{i,j}|^p] < \infty$. Then*

$$\mathbb{E}\left[\left\|\frac{1}{n}\sum_{i=1}^n X_i\right\|_p^p\right] \leq C_p \frac{\mathbb{E}\left[\|X_1\|_p^p\right]}{n^{p/2}}$$

*for every $p \geq 2$.*

*Proof of Corollary 2.* We know

$$\mathbb{E}\left[\left\|\frac{1}{n}\sum_{i=1}^{n}X_i\right\|_p^p\right] = \sum_{j=1}^{m}\mathbb{E}\left[\left|\frac{1}{n}\sum_{i=1}^{n}X_{i,j}\right|^p\right].$$

Applying Corollary 1 on each component of each $X_i$ yields

$$\sum_{j=1}^{m}\mathbb{E}\left[\left|\frac{1}{n}\sum_{i=1}^{n}X_{i,j}\right|^p\right] \leq C_p\sum_{j=1}^{m}\frac{\mathbb{E}|X_{1,j}|^p}{n^{p/2}} = C_p\frac{\mathbb{E}\left[\|X_1\|_p^p\right]}{n^{p/2}},$$

as desired. □

We also need the following inequality to compare $\|x\|_p$ and $\|x\|_q$ for $p \neq q$ and $x \in \mathbb{R}^m$. The proof follows directly from the Hölder's inequality.

**Lemma 3.** *For any $x \in \mathbb{R}^m$ and $p < q$, we have:*

$$\|x\|_p \leq m^{1/p-1/q}\|x\|_q.$$

*Proof.*

$$\|x\|_p^p = \sum_{i=1}^{m}|x_i|^p \cdot 1 \leq \left(\sum_{i=1}^{m}|x_i|^q\right)^{p/q} m^{1-p/q}$$

where the last inequality follows from the Hölder's inequality. Our result follows by taking the $(1/p)$-th power on both sides. □

## A.2  Bounding $\mathbb{E}[|\Delta_n|^2]$

Recall that $\Delta_n = g\left(S_H(2^n)/2^n\right) - \frac{1}{2}\left(g\left(S_H^O(2^{n-1})/2^{n-1}\right) + g\left(S_H^E(2^{n-1})/2^{n-1}\right)\right)$, and the final estimator takes the form $\Delta_N/p_N + g(H_1)$. Therefore, understanding the theoretical properties of $\Delta_n$ is crucial for studying our estimator.

*Proof of Lemma 1.* For simplicity, we denote $m(\pi)$ by $\mu$. By Assumption 3.3, there exists $\epsilon > 0$ such that $g$ is $\alpha$-Hölder continuous on $(\mu - \epsilon, \mu + \epsilon)$, we can then write $\Delta_n$ as:

$$|\Delta_n| = |\Delta_n|\mathbf{1}(A_1) + |\Delta_n|\mathbf{1}(A_2), \tag{15}$$

where $A_1$ is the event

$$\left\{\left\|\frac{S_H^O(2^{n-1})}{2^{n-1}} - \mu\right\| < \epsilon\right\} \cap \left\{\left\|\frac{S_H^E(2^{n-1})}{2^{n-1}} - \mu\right\| < \epsilon\right\},$$

and $A_2$ is the event

$$\left\{\max\left(\left\|\frac{S_H^O(2^{n-1})}{2^{n-1}} - \mu\right\|, \left\|\frac{S_H^E(2^{n-1})}{2^{n-1}} - \mu\right\|\right) \geq \epsilon\right\}.$$

Under the event $A_1$, we have $\left\| \frac{S_H^O(2^{n-1})}{2^{n-1}} - \mu \right\| < \epsilon$ and $\left\| \frac{S_H^E(2^{n-1})}{2^{n-1}} - \mu \right\| < \epsilon$. This further implies

$$\left\| \frac{S_H(2^n)}{2^n} - \mu \right\| < \epsilon$$

by the triangle inequality and the fact $\frac{S_H(2^n)}{2^n} = \frac{1}{2}\left( \frac{S_H^O(2^{n-1})}{2^{n-1}} + \frac{S_H^E(2^{n-1})}{2^{n-1}} \right)$.

Then we can write $\Delta_n$ as:

$$
\begin{aligned}
\Delta_n &= g\left( \frac{S_H(2^n)}{2^n} \right) - \frac{1}{2}\left( \frac{S_H^O(2^{n-1})}{2^{n-1}} + \frac{S_H^E(2^{n-1})}{2^{n-1}} \right) \\
&= \frac{1}{2}\left( g\left( \frac{S_H(2^n)}{2^n} \right) - g\left( \frac{S_H^O(2^{n-1})}{2^{n-1}} \right) \right) + \frac{1}{2}\left( g\left( \frac{S_H(2^n)}{2^n} \right) - g\left( \frac{S_H^E(2^{n-1})}{2^{n-1}} \right) \right) \\
&= \frac{1}{2} Dg(\xi_n^O)\left( \frac{S_H(2^n)}{2^n} - \frac{S_H^O(2^{n-1})}{2^{n-1}} \right) + \frac{1}{2} Dg(\xi_n^E)\left( \frac{S_H(2^n)}{2^n} - \frac{S_H^E(2^{n-1})}{2^{n-1}} \right) \\
&= \frac{1}{4}\left( Dg(\xi_n^O) - Dg(\xi_n^E) \right) \frac{S_H^E(2^{n-1}) - S_H^O(2^{n-1})}{2^{n-1}},
\end{aligned}
$$

where $\xi_n^O$ is a convex combination of $\frac{S_H(2^n)}{2^n}$ and $\frac{S_H^O(2^{n-1})}{2^{n-1}}$, $\xi_n^E$ is a convex combination of $\frac{S_H(2^n)}{2^n}$ and $\frac{S_H^E(2^{n-1})}{2^{n-1}}$ by the Multivariate Mean value Theorem. Under $A_1$, both $\xi_n^O$ and $\xi_n^E$ are within the $\epsilon$-neighbor of $\mu$, applying the $\alpha$-Hölder continuous assumption yields

$$
|\Delta_n| \leq c_1(\epsilon) \left\| \xi_n^O - \xi_n^E \right\|^\alpha \cdot \left\| \frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}} \right\| \leq c_2(\epsilon) \left\| \frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}} \right\|^{1+\alpha}.
$$

Then,

$$
\mathbb{E}\left[ |\Delta_n|^2 \mathbf{1}(A_1) \right] \leq c_2(\epsilon) \mathbb{E}\left[ \left\| \frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}} \right\|^{2(1+\alpha)} \right]. \tag{16}
$$

Since $S_H^O(2^{n-1})$ and $S_H^E(2^{n-1})$ are vectors in $\mathbb{R}^m$, applying Lemma 3 on $p = 2, q = 2(1+\alpha)$ gives:

$$
\left\| \frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}} \right\|^{2(1+\alpha)} \leq m^\alpha \left\| \frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}} \right\|_{2(1+\alpha)}^{2(1+\alpha)} \tag{17}
$$

Since $S_H^O(2^{n-1}) - S_H^E(2^{n-1})$ is the sum of $2^{n-1}$ i.i.d. random variables, each with the same distribution as $H_2 - H_1$, applying the Multivariate Marcinkiewicz-Zygmund inequality (Corollary 2) gives us:

$$
\mathbb{E}\left[ \left\| \frac{S_H^O(2^{n-1}) - S_H^E(2^{n-1})}{2^{n-1}} \right\|_{2(1+\alpha)}^{2(1+\alpha)} \right] \leq C_{2(1+\alpha)} \cdot \frac{\mathbb{E}\left[ \|H_2 - H_1\|_{2(1+\alpha)}^{2(1+\alpha)} \right]}{2^{(1+\alpha)(n-1)}} \tag{18}
$$

$$
\leq C_{2(1+\alpha)} \cdot 2^{3(1+\alpha)} \cdot \frac{\mathbb{E}\left[ \|H_1\|_{2(1+\alpha)}^{2(1+\alpha)} \right]}{2^{(1+\alpha)n}}. \tag{19}
$$

where the last step uses the inequality $(a + b)^p \leq 2^{p-1}(|a|^p + |b|^p)$ for $p \geq 2$. It is worth mentioning that the right hand side of (19) is finite as Assumption 3.4 guarantees $H_1$ has finite $l$-th moment with $l > 2 + \alpha$. Combining (16), (17), and (19), we have:

$$\mathbb{E}\left[\|\Delta_n\|^2 \mathbf{1}(A_1)\right] \leq C_1(m, \alpha, \epsilon)2^{-n(1+\alpha)}, \tag{20}$$

where $C_1(m, \alpha, \epsilon) = c_2(\epsilon) \cdot C_{2(1+\alpha)} \cdot 2^{3(1+\alpha)} \cdot \mathbb{E}\left[\|H_1\|_{2(1+\alpha)}^{2(1+\alpha)}\right]$ is a constant when Assumption 3.1 − 3.4 are satisfied.

Under $A_2$, we have:

$$|\Delta_n|^2 \mathbf{1}(A_2) \leq |\Delta_n|^2 \mathbf{1}\left(\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right) + |\Delta_n|^2 \mathbf{1}\left(\left\|\frac{S_H^{\mathsf{E}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right) \tag{21}$$

Now we upper bound the first term's expectation,

$$\mathbb{E}\left[|\Delta_n|^2 \mathbf{1}\left(\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right)\right] \leq \mathbb{E}[|\Delta_n|^{2s}]^{1/s}\mathbb{P}\left(\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right)^{(s-1)/s} \tag{22}$$

$$\leq \mathcal{C}_s^{1/s}2^{-\alpha_s n/s}\mathbb{P}\left(\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right)^{(s-1)/s} \tag{23}$$

$$\leq \mathcal{C}_s^{1/s} \cdot (\epsilon^{-l(s-1)/s}) \cdot 2^{-\alpha_s n/s} \cdot \mathbb{E}\left[\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\|^l\right]^{(s-1)/s} . \tag{24}$$

Here (22) follows from the Hölder's inequality, (23) uses Assumption 3.5, and (24) follows from the Markov's inequality. Again, using Lemma 3 and Corollary 2, the term $\mathbb{E}\left[\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\|^l\right]$ can be upper bounded by:

$$\mathbb{E}\left[\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\|^l\right] \leq m^{l/2-1}\mathbb{E}\left[\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\|_l^l\right] \leq 2^{l/2} \cdot m^{l/2-1} \cdot C_l \cdot \frac{\mathbb{E}\left[\|H_1\|_l^l\right]}{2^{nl/2}}. \tag{25}$$

Combining (24) and (25), we have

$$\mathbb{E}\left[|\Delta_n|^2 \mathbf{1}\left(\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right)\right] \leq C_2(m, l, \epsilon, s)2^{-\alpha_s n/s}2^{-nl(s-1)/(2s)}$$

$$= C_2(m, l, \epsilon, s)2^{-n\left(\frac{\alpha_s}{s} + \frac{(s-1)l}{2s}\right)},$$

where $C_2(m, l, \epsilon, s) = \mathcal{C}_s^{1/s} \cdot \left(\epsilon^{-l}2^{l/2} \cdot m^{l/2-1} \cdot C_l \cdot \mathbb{E}\left[\|H_1\|_l^l\right]\right)^{(s-1)/s}$ is a constant when Assumption 3.1 − 3.5 are satisfied. Furthermore, by Assumption 3.5, $2\alpha_s + (s-1)l > 2s$. It is clear that $\frac{\alpha_s}{s} + \frac{(s-1)l}{2s} > 1$, and therefore

$$\mathbb{E}\left[|\Delta_n|^2 \mathbf{1}\left(\left\|\frac{S_H^{\mathsf{O}}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right)\right] \leq C_2(m, l, \epsilon, s)2^{-(1+\tilde{\alpha})n}, \tag{26}$$

where $\tilde{\alpha} = \frac{\alpha_s}{s} + \frac{(s-1)l}{2s} - 1 > 0$. The same argument also shows

$$\mathbb{E}\left[|\Delta_n|^2 \mathbf{1}\left(\left\|\frac{S_H^\mathbf{O}(2^{n-1})}{2^{n-1}} - \mu\right\| > \epsilon\right)\right] \leq C_2(m, l, \epsilon, s) 2^{-(1+\tilde{\alpha})n}. \tag{27}$$

Combining (26), (27), and (21), we have

$$\mathbb{E}\left[|\Delta_n|^2 \mathbf{1}(A_2)\right] \leq 2C_2(m, l, \epsilon, s) 2^{-(1+\tilde{\alpha})n}. \tag{28}$$

Finally, taking $\gamma = \min\{\alpha, \tilde{\alpha}\}$, $C = C_1 + 2C_2$, and using (15), (20), and (28), we conclude:

$$\mathbb{E}[|\Delta_n|^2] \leq C 2^{-n(1+\gamma)}. \tag{29}$$

$\square$

## A.3 The Moment Assumption 3.4 and Markov chain mixing

In this subsection we discuss the relation between the Moment Assumption 3.4 and the mixing time of the underlying Markov chain. Throughout this subsection, the unbiased estimator $H$ of $m(\pi)$ is assumed to be the JOA estimator $H_k(Y, Z)$ defined in (3), which also extends to $H_{k:m}(Y, Z) = (m - k + 1)^{-1} \sum_{l=k}^{m} H_l(Y, Z)$ naturally.

Before giving a formal statement of Proposition 3, we first recall some definitions in Markov chain theory. We say a $\pi$-invariant, $\phi$-irreducible and aperiodic Markov transition kernel $P$ satisfies a geometric drift condition if there exists a measurable function $V : \Omega \to [1, \infty)$, $\lambda \in (0, 1)$, and a measurable set $\mathcal{S}$ such that for all $x \in \Omega$:

$$\int P(x, \mathrm{d}y)V(y) \leq \lambda V(x) + b\mathbf{1}(x \in \mathcal{S}). \tag{30}$$

Moreover, the set $\mathcal{S}$ is called a small set if there exists a positive integer $m$, $\epsilon > 0$, and a probability measure $\nu$ on such that for every $x \in \mathcal{S}$:

$$P^m(x, \cdot) \geq \epsilon \mu(\cdot). \tag{31}$$

The technical definitions for irreducibility, aperiodicity and small sets can be found in Chapter 5 of Meyn and Tweedie [2012]. The geometric drift condition is a key tool guaranteeing the geometric ergodicity of a Markov chain, meaning the Markov chain $P$ converges to its stationary distribution $\pi$ at a geometric rate (see Theorem 9 of Roberts and Rosenthal [2004], Jones and Hobert [2001]). It is known that the geometric drift condition is satisfied for a wide family of Metropolis-Hastings algorithms, we refer the readers to Mengersen and Tweedie [1996], Roberts and Tweedie [1996a,b], Jarner and Hansen [2000], and Roberts and Rosenthal [1997] for existing results.

Now we give a formal statement of Proposition 3.

**Proposition 4** (Verifying Assumption 3.4, formal version of Proposition 3). *Suppose the Markov transition kernel described in Section 3.1.1 satisfies a geometric drift condition with a small set $\mathcal{S}$ of the form $\mathcal{S} = \{x : V(x) \leq L\}$ for $\lambda + b/(1 + L) < 1$. Suppose there exists $\tilde{\epsilon} \in (0, 1)$ such that*

$$\inf_{(x,y) \in \mathcal{S} \times \mathcal{S}} \bar{P}((x, y), \mathcal{D}) \geq \tilde{\epsilon},$$

*where $\mathcal{D} := \{(x, x) : x \in \Omega\}$ is the diagonal of $\Omega \times \Omega$. Suppose also there exists $p > l$ and $D_p > 0$ such that $\mathbb{E}[\|f(Y_t)\|_p^p] < D_p$ for every $t$. Then $\mathbb{E}[\|H_k(Y, Z)\|_l^l] < \infty$ for every $k$.*

The main ingredient in the proof of Proposition 4 is to control the tail probability of the meeting time $\tau$. We say $\tau$ has a $\beta$-polynomial tail if there exists a constant $K_\beta > 0$ such that

$$\mathbb{P}(\tau > n) \leq K_\beta n^{-\beta}. \tag{32}$$

We say $\tau$ has an exponential tail if there exists a constant $K > 0$ and $\gamma \in (0, 1)$ such that

$$\mathbb{P}(\tau > n) \leq K\gamma^n. \tag{33}$$

Our next result gives sufficient conditions to ensure Assumption 3.4.

**Lemma 4.** *Suppose one of the following holds:*

- *There exists $p > l$, $\beta > 0$, and $D_p > 0$ such that $\frac{1}{p} + \beta > \frac{1}{l}$; $\mathbb{E}[\|f(Y_t)\|_p^p] < D_p$ for every $t$, and $\tau$ has a $\beta$-polynomial tail;*
- *There exists $p > l$ and $D_p > 0$ such that $\mathbb{E}[\|f(Y_t)\|_p^p] < D_p$ for every $t$, and $\tau$ has an exponential tail.*

*Then $\mathbb{E}[\|H_k(Y, Z)\|_l^l] < \infty$ for every $k$.*

*Proof of Lemma 4.* We start with the first case. Without loss of generality, we assume $k = 0$ and the estimator $H_0(Y, Z) := f(Y_0) + \sum_{i=1}^{\tau-1}(f(Y_i) - f(Z_{i-1}))$ takes scalar value. Let $D_k := f(Y_k) - f(Z_{k-1})$ for $k \geq 1$, and $D_0 = f(Y_0)$, the estimator can be written as:

$$H_0(Y, Z) = \sum_{k=0}^{\infty} D_k \mathbf{1}(\tau > k).$$

The meeting time $\tau$ is almost surely (a.s.) finite by the $\beta$-polynomial assumption, therefore $H_0(Y, Z)$ is the limit of $H_0^n(Y, Z) := \sum_{k=0}^{n} D_k \mathbf{1}(\tau > k)$ in the a.s. sense. We will now prove $H_0^n(Y, Z) \to H_0(Y, Z)$ in $L^l$, which further implies $\mathbb{E}[|H_0(Y, Z)|^l] < \infty$.

By the Minkowski's inequality on the probability space $L^l(\Omega)$, we have

$$\left(\mathbb{E}[|H_0^n(Y, Z) - H_0(Y, Z)|^l]\right)^{1/l} = \left(\mathbb{E}[|\sum_{k=n+1}^{\infty} D_k \mathbf{1}(\tau > k)|^l]\right)^{1/l} \tag{34}$$

$$\leq \sum_{k=n+1}^{\infty} \left(\mathbb{E}[|D_k \mathbf{1}(\tau > k)|^l]\right)^{1/l}. \tag{35}$$

Every term in (35) can be upper bounded by the Hölder's inequality

$$\left(\mathbb{E}[|D_k \mathbf{1}(\tau > k)|^l]\right)^{1/l} \leq \left(\mathbb{E}[|D_k|^p]\right)^{1/p}\left(\mathbb{P}(\tau > k)\right)^{1/q} \quad \text{here } 1/q = 1/l - 1/p \tag{36}$$

$$\leq (2D_p)^{1/p} K_\beta^{1/q} k^{-\beta/q} \tag{37}$$

$$= (2D_p)^{1/p} K_\beta^{1/q} k^{-\frac{\beta}{\frac{1}{l} - \frac{1}{p}}}. \tag{38}$$

Since $\beta > \frac{1}{l} - \frac{1}{p} > 0$, the right hand side of (37) is summable. Therefore we conclude

$$\sum_{k=n+1}^{\infty} \left(\mathbb{E}[|D_k \mathbf{1}(\tau > k)|^l]\right)^{1/l} \to 0$$

33

as $n \to \infty$, so $H_0^n(Y, Z) \to H_0(Y, Z)$ in $L^l$.

In the second case, exponential light tail implies $\beta$-polynomial tail for every $\beta > 0$, our result immediately follows from the first case.

$\square$

The assumption $\mathbb{E}[\|f(Y_t)\|^p] < D_p$ in Lemma 4 is generally satisfied as long as $f$ has $p$-th moment under the stationary distribution $\pi$. It remains to verify the tail conditions of $\tau$, i.e., formula (32) or (33). The exponential tail (33) and polynomial tail (32) are closely related to the geometric ergodicity and polynomial ergodicity of the underlying marginal Markov chain $P$, respectively. For simplicity, we only give conditions for the exponential tail here, which is provided in Jacob et al. [2020]. The sufficient conditions of polynomial tail of $\tau$ can be founded in Theorem 2 of Middleton et al. [2020].

**Proposition 5** (Proposition 3.4 in Jacob et al. [2020])**.** *Suppose the Markov transition kernel described in Section 3.1.1 satisfies a geometric drift condition with a small set $\mathcal{S}$ of the form $\mathcal{S} = \{x : V(x) \leq L\}$ for $\lambda + b/(1 + L) < 1$. Suppose there exists $\tilde{\epsilon} \in (0, 1)$ such that*

$$\inf_{(x,y) \in \mathcal{S} \times \mathcal{S}} \bar{P}((x, y), \mathcal{D}) \geq \tilde{\epsilon},$$

*where $\mathcal{D} := \{(x, x) : x \in \Omega\}$ is the diagonal of $\Omega \times \Omega$. Then the meeting time $\tau$ has a exponential light tail.*

Combining Lemma 4 and Proposition 5, the proof of Proposition 4 is immediate.

*Proof of Proposition 4.* By Proposition 5, we know $\tau$ has an exponential tail. Using the second case of Lemma 4, our result follows. $\square$

It is still possible to further strengthen Proposition 4 given extra assumptions on $\tau$ or $f$. For example, when $\tau$ has an exponential tail and $\mathbb{E}_\pi[e^{\theta f}] < \infty$ for a univariate $f$ and some $\theta > 0$, one can then prove the JOA estimator also has an exponential moment, and thus has every finite-order moment. The existence of an exponential moment may be helpful for analyzing the concentration properties of the JOA estimator.

# B  Delta method for Section 4.2.2

To give a confidence interval for $Z(\theta_2)/Z(\theta_1)$ by Gibbs Sampling, we first write $\pi_{\theta_1} = \mathbb{E}_{\theta_1}[e^{\theta_1 H(\sigma)}], \pi_{\theta_2} = \mathbb{E}_{\theta_2}[e^{\theta_2 H(\sigma)}]$, and $g(x, y) = y/x$. Then the quantity of our interest can be written as

$$\frac{Z(\theta_2)}{Z(\theta_1)} = g(\pi_{\theta_1}, \pi_{\theta_2}).$$

The gradient of $g$ can be calculated by:

$$\nabla g = \left(-y/x^2, \ 1/x\right)^\top.$$

Let $\hat{\pi}_{\theta_1}, \hat{\pi}_{\theta_2}$ be the empirical averages of the Gibbs samplers, $\sigma_{\theta_1}^2, \sigma_{\theta_2}^2$ their asymptotic variances, and $\hat{\sigma}_{\theta_1}^2, \hat{\sigma}_{\theta_2}^2$ their empirical variances. Then the delta method implies:

$$\sqrt{n}g(\hat{\pi}_{\theta_1}, \hat{\pi}_{\theta_2}) - g(\pi_{\theta_1}, \pi_{\theta_2})) \to \mathsf{N}\left(0, \nabla g(\pi_{\theta_1}, \pi_{\theta_2})^\top \begin{bmatrix} \sigma_{\theta_1}^2 & 0 \\ 0 & \sigma_{\theta_2}^2 \end{bmatrix} \nabla g(\pi_{\theta_1}, \pi_{\theta_2})\right)$$

provided the Gibbs sampler satisfies a Markov Chain Central Limit Theorem [Jones, 2004].

Therefore, let

$$\hat{\sigma}_R^2 := \nabla g(\hat{\pi}_{\theta_1}, \hat{\pi}_{\theta_2})^\top \begin{bmatrix} \hat{\sigma}_{\theta_1}^2 & 0 \\ 0 & \hat{\sigma}_{\theta_2}^2 \end{bmatrix} \nabla g(\hat{\pi}_{\theta_1}, \hat{\pi}_{\theta_2})$$

be the estimator for the variance term. The $100(1-\alpha)\%$ Confidence interval of the Gibbs sampler can be constructed by

$$\left(\frac{\hat{\pi}_{\theta_2}}{\hat{\pi}_{\theta_1}} - Z_{\alpha/2} \cdot \hat{\sigma}_R, \frac{\hat{\pi}_{\theta_2}}{\hat{\pi}_{\theta_1}} + Z_{\alpha/2} \cdot \hat{\sigma}_R\right),$$

where $Z_{\alpha/2}$ is the $(\alpha/2)$-th quantile of the standard Normal distribution.

# References

A. Ades, G. Lu, and K. Claxton. Expected value of sample information calculations in medical decision modeling. *Medical decision making*, 24(2):207–227, 2004. 11

S. Agapiou, G. O. Roberts, and S. J. Vollmer. Unbiased Monte Carlo: Posterior estimation for intractable/infinite-dimensional models. *Bernoulli*, 24(3):1726–1786, 2018. 9

C. Andrieu and G. O. Roberts. The pseudo-marginal approach for efficient Monte Carlo computations. *The Annals of Statistics*, 37(2):697–725, 2009. 3

D. Belomestny, M. Ladkau, and J. Schoenmakers. Multilevel simulation based policy iteration for optimal stopping–convergence and complexity. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):460–483, 2015. 2, 4, 11

C. H. Bennett. Efficient estimation of free energy differences from Monte Carlo data. *Journal of Computational Physics*, 22(2):245–268, 1976. 21

A. Beskos, A. Jasra, K. Law, R. Tempone, and Y. Zhou. Multilevel sequential Monte Carlo samplers. *Stochastic Processes and their Applications*, 127(5):1417–1440, 2017. 4

N. Biswas and L. Mackey. Bounding Wasserstein distance with couplings. *arXiv preprint arXiv:2112.03152*, 2021. 3

N. Biswas, P. E. Jacob, and P. Vanetti. Estimating convergence of Markov chains with L-lag couplings. In *Advances in Neural Information Processing Systems*, volume 32, 2019. 3, 8

J. H. Blanchet and P. W. Glynn. Unbiased Monte Carlo for optimization and functions of expectations via multi-level randomization. *2015 Winter Simulation Conference (WSC)*, pages 3656–3667, 2015. 4, 5, 9, 14, 15, 18

J. H. Blanchet, N. Chen, and P. W. Glynn. Unbiased monte carlo computation of smooth functions of expectations via taylor expansions. In *2015 Winter Simulation Conference (WSC)*, pages 360–367. IEEE, 2015. 3, 5, 9

J. H. Blanchet, P. W. Glynn, and Y. Pei. Unbiased Multilevel Monte Carlo: Stochastic optimization, steady-state simulation, quantiles, and other applications. *arXiv preprint arXiv:1904.09929*, 2019. 4, 15

W. Chen, L. Ma, and X. Liang. Blind identification based on expectation-maximization algorithm coupled with blocked Rhee–Glynn smoothing estimator. *IEEE Communications Letters*, 22(9):1838–1841, 2018. 3

P. Diaconis and G. Wang. Bayesian goodness of fit tests: a conversation for David Mumford. *Annals of Mathematical Sciences and Applications*, 3(1):287–308, 2018. 21

C. R. Doss, J. M. Flegal, G. L. Jones, and R. C. Neath. Markov chain Monte Carlo estimation of quantiles. *Electronic Journal of Statistics*, 8(2):2448–2478, 2014. 2

R. Douc, P. E. Jacob, A. Lee, and D. Vats. Unbiased estimation of the asymptotic variance of MCMC estimators using coupled chains. *forthcoming*, 2021. 3

R. Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019. 17

P. Fearnhead, O. Papaspiliopoulos, G. O. Roberts, and A. Stuart. Random-weight particle filtering of continuous time processes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):497–512, 2010. 9

M. B. Giles. Multilevel Monte Carlo path simulation. *Operations Research*, 56(3):607–617, 2008. 3, 4

M. B. Giles. Multilevel Monte Carlo methods. *Acta Numer.*, 24:259–328, 2015. 3

M. B. Giles. MLMC for nested expectations. In *Contemporary Computational Mathematics-A Celebration of the 80th Birthday of Ian Sloan*, pages 425–442. Springer, 2018. 11

M. B. Giles and T. Goda. Decision-making under uncertainty: using MLMC for efficient estimation of EVPPI. *Statistics and Computing*, 29(4):739–751, 2019. 2, 11, 24

P. W. Glynn and P. Heidelberger. Analysis of parallel replicated simulations under a completion time constraint. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 1(1):3–23, 1991. 18

P. W. Glynn and C.-h. Rhee. Exact estimation for Markov chain equilibrium expectations. *Journal of Applied Probability*, 51(A):377–389, 2014. 3, 9

P. W. Glynn and W. Whitt. The asymptotic efficiency of simulation estimators. *Operations research*, 40(3):505–520, 1992. 18

T. Goda, T. Hironaka, W. Kitade, and A. Foster. Unbiased MLMC stochastic gradient-based optimization of Bayesian experimental designs. *SIAM Journal on Scientific Computing*, 44(1):A286–A311, 2022. 11

S. Heinrich. Multilevel Monte Carlo methods. In *International Conference on Large-Scale Scientific Computing*, pages 58–67. Springer, 2001. 3

J. Heng and P. E. Jacob. Unbiased Hamiltonian Monte Carlo with couplings. *Biometrika*, 106(2):287–302, 2019. 4

T. Hironaka and T. Goda. An efficient estimation of nested expectations without conditional sampling. *arXiv preprint arXiv:2111.12278*, 2021. 11

P. E. Jacob and A. H. Thiery. On nonnegative unbiased estimators. *The Annals of Statistics*, 43(2):769–784, 2015. 5, 12, 27

P. E. Jacob, J. O'Leary, and Y. F. Atchadé. Unbiased markov chain monte carlo methods with couplings. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(3):543–600, 2020. 3, 4, 7, 9, 14, 17, 19, 20, 21, 25, 26, 34

S. F. Jarner and E. Hansen. Geometric ergodicity of Metropolis algorithms. *Stochastic processes and their applications*, 85(2):341–361, 2000. 14, 32

A. Jasra, K. Kamatani, K. J. Law, and Y. Zhou. A multi-index Markov chain Monte Carlo method. *International Journal for Uncertainty Quantification*, 8(1), 2018. 4

Y. H. Jiang, T. Liu, Z. Lou, J. S. Rosenthal, S. Shangguan, F. Wang, and Z. Wu. Mcmc Confidence Intervals and Biases. *arXiv preprint arXiv:2012.02816*, 2020. 23

G. L. Jones. On the Markov chain central limit theorem. *Probability surveys*, 1:299–320, 2004. 23, 35

G. L. Jones and J. P. Hobert. Honest exploration of intractable probability distributions via Markov chain Monte Carlo. *Statistical Science*, pages 312–334, 2001. 32

M. Keane and G. L. O'Brien. A bernoulli factory. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 4(2):213–219, 1994. 5

R. Koenker and K. F. Hallock. Quantile regression. *Journal of economic perspectives*, 15 (4):143–156, 2001. 2

D. Levy, Y. Carmon, J. C. Duchi, and A. Sidford. Large-Scale Methods for Distributionally Robust Optimization. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*, 2020. 3

S. Livingstone, M. Betancourt, S. Byrne, and M. Girolami. On the geometric ergodicity of Hamiltonian Monte Carlo. *Bernoulli*, 25(4A):3109–3138, 2019. 14

A.-M. Lyne, M. Girolami, Y. Atchadé, H. Strathmann, and D. Simpson. On russian roulette estimates for bayesian inference with doubly-intractable likelihoods. *Statistical science*, 30(4):443–467, 2015. 11, 27

J. Marcinkiewicz and A. Zygmund. Quelques théoremes sur les fonctions indépendantes. *Fund. Math*, 29:60–90, 1937. 28

D. Maucort-Boulch, S. Franceschi, and M. Plummer. International correlation between human papillomavirus prevalence and cervical cancer incidence. *Cancer Epidemiology and Prevention Biomarkers*, 17(3):717–720, 2008. 25

D. McLeish. A general method for debiasing a Monte Carlo estimator. *Monte Carlo Methods and Applications*, 17(4):301–315, 2011. 3

X.-L. Meng and W. H. Wong. Simulating ratios of normalizing constants via a simple identity: a theoretical exploration. *Statistica Sinica*, pages 831–860, 1996. 2, 11, 21

K. L. Mengersen and R. L. Tweedie. Rates of convergence of the Hastings and Metropolis algorithms. *The Annals of Statistics*, 24(1):101–121, 1996. 14, 32

S. P. Meyn and R. L. Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012. 32

L. Middleton, G. Deligiannidis, A. Doucet, and P. E. Jacob. Unbiased Markov chain Monte Carlo for intractable target distributions. *Electronic Journal of Statistics*, 14 (2):2842–2891, 2020. 4, 7, 34

I. Murray, Z. Ghahramani, and D. J. MacKay. MCMC for doubly-intractable distributions. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, pages 359–366, 2006. 11

Ş. Nacu and Y. Peres. Fast simulation of new coins from old. *The Annals of Applied Probability*, 15(1A):93–115, 2005. 5

J. O'Leary and G. Wang. Metropolis-Hastings transition kernel couplings. *arXiv preprint arXiv:2102.00366*, 2021. 8

O. Papaspiliopoulos. A methodological framework for Monte Carlo probabilistic inference for diffusion processes. 2009. 9

M. Plummer. Cuts in Bayesian graphical models. *Statistics and Computing*, 25(1):37–43, 2015. 11, 25

E. Pompe and P. E. Jacob. Asymptotics of cut distributions and robust modular inference using Posterior Bootstrap. *arXiv preprint arXiv:2110.11149*, 2021. 11, 25

J. G. Propp and D. B. Wilson. Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures & Algorithms*, 9(1-2):223–252, 1996. 21

T. Rainforth, R. Cornish, H. Yang, A. Warrington, and F. Wood. On nesting Monte Carlo estimators. In *International Conference on Machine Learning*, pages 4267–4276. PMLR, 2018. 2, 11

C.-h. Rhee and P. W. Glynn. Unbiased estimation with square root convergence for SDE models. *Operations Research*, 63(5):1026–1043, 2015. 3, 4, 15

M. Rischard, P. E. Jacob, and N. Pillai. Unbiased estimation of log normalizing constants with applications to Bayesian cross-validation. *arXiv preprint arXiv:1810.01382*, 2018. 3

G. Roberts and J. Rosenthal. Geometric ergodicity and hybrid Markov chains. *Electronic Communications in Probability*, 2:13–25, 1997. 32

G. O. Roberts and J. S. Rosenthal. General state space Markov chains and MCMC algorithms. *Probability surveys*, 1:20–71, 2004. 32

G. O. Roberts and R. L. Tweedie. Exponential convergence of langevin distributions and their discrete approximations. *Bernoulli*, pages 341–363, 1996a. 14, 32

G. O. Roberts and R. L. Tweedie. Geometric convergence and central limit theorems for multidimensional Hastings and Metropolis algorithms. *Biometrika*, 83(1):95–110, 1996b. 14, 32

Y. Romano, E. Patterson, and E. Candes. Conformalized quantile regression. *Advances in Neural Information Processing Systems*, 32:3543–3553, 2019. 2

J. S. Rosenthal. Faithful couplings of Markov chains: now equals forever. *Advances in Applied Mathematics*, 18(3):372–381, 1997. 7

J. S. Rosenthal. Simple confidence intervals for MCMC without CLTs. *Electronic Journal of Statistics*, 11(1):211–214, 2017. 23

F. J. Ruiz, M. K. Titsias, T. Cemgil, and A. Doucet. Unbiased gradient estimation for variational auto-encoders using coupled markov chains. *arXiv preprint arXiv:2010.01845*, 2020. 3

Y. Shi and R. Cornish. On Multilevel Monte Carlo Unbiased Gradient Estimation for Deep Latent Variable Models. In *International Conference on Artificial Intelligence and Statistics*, pages 3925–3933. PMLR, 2021. 3

V. B. Tadić and A. Doucet. Asymptotic bias of stochastic gradient search. *The Annals of Applied Probability*, 27(6):3255–3304, 2017. 3

I. Takeuchi, Q. Le, T. Sears, and A. Smola. Nonparametric quantile estimation. 2006. 2

M. Vihola. Unbiased estimators and multilevel Monte Carlo. *Operations Research*, 66(2): 448–462, 2018. 5, 14, 15

W. Wagner. Unbiased Monte Carlo evaluation of certain functional integrals. *Journal of Computational Physics*, 71(1):21–33, 1987. 9

G. Wang. On the theoretical properties of the exchange algorithm. *arXiv preprint arXiv:2005.09235*, 2020. 14

G. Wang, J. O'Leary, and P. Jacob. Maximal Couplings of the Metropolis-Hastings Algorithm. In *International Conference on Artificial Intelligence and Statistics*, pages 1225–1233. PMLR, 2021. 8

Z. Zhou, G. Wang, J. Blanchet, and P. W. Glynn. Unbiased Optimal Stopping via the MUSE. *arXiv preprint arXiv:2106.02263*, 2021. 2, 3, 4, 11, 14