

## Variance Reduction for Dependent Sequences with Applications to Stochastic Gradient MCMC\*

Denis Belomestny<sup>†</sup>, Leonid Iosipoi<sup>‡</sup>, Eric Moulines<sup>§</sup>, Alexey Naumov<sup>‡</sup>, and  
Sergey Samsonov<sup>‡</sup>

**Abstract.** In this paper we propose a novel and practical variance reduction approach for additive functionals of dependent sequences. Our approach combines the use of control variates with the minimization of an empirical variance estimate. We analyze finite sample properties of the proposed method and derive finite-time bounds of the excess asymptotic variance to zero. We apply our methodology to stochastic gradient Markov chain Monte Carlo (SGMCMC) methods for Bayesian inference on large data sets and combine it with existing variance reduction methods for SGMCMC. We present empirical results carried out on a number of benchmark examples showing that our variance reduction method achieves significant improvement as compared to state-of-the-art methods at the expense of a moderate increase of computational overhead.

**Key words.** MCMC algorithms, variance reduction, stochastic gradient, generative models

**AMS subject classifications.** 60J20, 65C40, 65C60

**DOI.** 10.1137/19M1301199

**1. Introduction.** Variance reduction aims at reducing the stochastic error of a Monte Carlo estimate; see [39], [42], [25], and [24] for an introduction to this field. Recently one witnessed a revival of interest in variance reduction techniques for dependent sequences with applications to Bayesian inference and reinforcement learning among others; see, for instance, [33], [28], [15], [11], [2], and references therein.

Suppose that we wish to compute the integral of an arbitrary function  $f : \mathcal{X} \mapsto \mathbb{R}$  with respect to a probability measure  $\pi$  on a general state-space  $(\mathcal{X}, \mathcal{X})$ , that is,  $\pi(f) = \int_{\mathcal{X}} f(x)\pi(dx)$ . If drawing an independent and identically distributed (i.i.d.) sample from  $\pi$  is an option, a natural estimator for  $\pi(f)$  is the sample mean

$$\pi_N(f) := N^{-1} \sum_{k=0}^{N-1} f(X_k), \quad N \in \mathbb{N},$$

where  $(X_k)_{k=0}^{N-1}$  is an i.i.d. sample from  $\pi$ . Using the central limit theorem, one can construct an asymptotically valid confidence interval for the value  $\pi(f)$  of the form  $\pi_N(f) \pm$

\*Received by the editors November 20, 2019; accepted for publication (in revised form) January 22, 2021; published electronically May 3, 2021.

<https://doi.org/10.1137/19M1301199>

**Funding:** This work was prepared within the framework of the HSE University Basic Research Program. Results of section 3 were obtained under support of the RSF grant 19-71-30020 (HSE University).

<sup>†</sup>Duisburg-Essen University, Essen, Germany, and HSE University, Moscow, Russia ([denis.belomestny@uni-due.de](mailto:denis.belomestny@uni-due.de)).

<sup>‡</sup>HSE University, Moscow, Russia ([liosipoi@hse.ru](mailto:liosipoi@hse.ru), [anaumov@hse.ru](mailto:anaumov@hse.ru), [svsamsonov@hse.ru](mailto:svsamsonov@hse.ru)).

<sup>§</sup>Ecole Polytechnique, Paris, France, and HSE University, Moscow, Russia ([eric.moulines@polytechnique.edu](mailto:eric.moulines@polytechnique.edu)).

$\mathfrak{q} N^{-1/2}(\text{Var}_\pi(f))^{1/2}$ , where  $\mathfrak{q}$  is a quantile of a normal distribution, and  $\text{Var}_\pi(f) = \int_{\mathbf{X}} \{f(x) - \pi(f)\}^2 \pi(dx)$ . A general way to reduce the variance  $\text{Var}_\pi(f)$  is to select another function  $g$  in a set  $\mathcal{G}$  such that  $\pi(g) = 0$  and  $\text{Var}_\pi(f - g) \ll \text{Var}_\pi(f)$ . Such a function  $g$  is called a *control variate*. A natural approach to learn  $g \in \mathcal{G}$  is to minimize the empirical variance

$$(1.1) \quad D_n(f - g) = (n - 1)^{-1} \sum_{k=0}^{n-1} (f(X_k) - g(X_k) - \pi_n(f - g))^2$$

constructed using a new independent learning sample  $(X_k)_{k=0}^{n-1}$ . This leads to the empirical variance minimization (EVM) method recently studied in [6] and [7]. In many problems of interest, drawing an i.i.d. sample from  $\pi$  is not an option, yet it is possible to obtain a nonstationary dependent sequence  $(X_k)_{k=0}^\infty$  whose marginal distribution converges to  $\pi$ . This situation is typical in Bayesian statistics, where  $\pi$  represents a posterior distribution and  $(X_k)_{k=0}^\infty$  is sampled using Markov chain Monte Carlo (MCMC) methods. Under appropriate conditions, the central limit theorem also holds, and therefore, it is possible to construct the asymptotic confidence interval for  $\pi(f)$  of the form

$$(1.2) \quad \left[ \pi_N(f) - \mathfrak{q} \sqrt{\frac{V_\infty(f)}{N}}, \pi_N(f) + \mathfrak{q} \sqrt{\frac{V_\infty(f)}{N}} \right],$$

where  $V_\infty(f)$  is the asymptotic variance defined as

$$(1.3) \quad V_\infty(f) := \lim_{N \rightarrow \infty} N \cdot \mathbb{E} \left[ (\pi_N(f) - \pi(f))^2 \right].$$

A sensible approach is to select a control variate  $g \in \mathcal{G}$  by minimizing an estimate for the asymptotic variance  $V_\infty(f - g)$ . When the spectral estimate of  $V_\infty(f - g)$  is used, this leads to the empirical spectral variance minimization (ESVM); see [5].

In this paper, special attention is paid to the case when  $\mathbf{X} = \mathbb{R}^d$  and  $\pi$  admits a smooth and everywhere positive density (also denoted by  $\pi$ ) w.r.t. the Lebesgue measure such that the gradient  $\nabla U := -\nabla \log \pi$  can be evaluated. We study below sampling methods derived from the discretization of the overdamped Langevin dynamics. It is defined by the following stochastic differential equation:

$$(1.4) \quad dY_t = -\nabla U(Y_t) dt + \sqrt{2} dW_t,$$

where  $(W_t)_{t \geq 0}$  is the standard Brownian motion. Note that  $\nabla U$  does not depend on the normalizing constant of  $\pi$  which is typically unknown in Bayesian inference. Under some technical conditions, the distribution of  $Y_t$  converges to  $\pi$  as  $t \rightarrow \infty$ ; see [40]. The gradient-based MCMC algorithms are based on a time-discretized version of (1.4). In the Bayesian setting, a computational bottleneck of these algorithms is that the complexity of the gradient  $\nabla U$  evaluation scales proportionally to the number of observations (sample size)  $K$  which can be very time consuming in the “big data” limit. To alleviate this problem, [46] proposed to replace the “full” gradient  $\nabla U$  by a stochastic gradient estimate based on sums over random *minibatches*. This algorithm, stochastic gradient Langevin dynamics (SGLD), has emerged as a key MCMC algorithm in Bayesian inference for large scale datasets. The analysis of SGLD and its finite sample performance has attracted a wealth of contributions; see, for example, [30], [45], [34], [14], and the references therein. These works show that the use of stochastic

gradient comes at a price: while the resulting estimate of the gradient is still unbiased, its variance might annihilate the computational advantages of SGLD [14]. Several proposals have been made to reduce the variance of the stochastic gradient estimate of the “full” gradient, inspired by several methods, proposed for incremental stochastic optimization; see [41], [28], and [15]. [20] has investigated the properties of the Stochastic Average Gradient (SAGA) and Stochastic Variance Reduced Gradient (SVRG) estimators for Langevin dynamics. These results have been later completed and sharpened by [14], [11], [10]. Other variance reduction approaches include various subsampling schemes and constructing alternative estimates for the gradient (see, for instance, [2] and [47]).

The paper is organized as follows. In section 2, we analyze the ESVM approach for general dependent sequences. In particular, the ESVM method is described in subsection 2.1. In subsection 2.2, we study the theoretical properties of the ESVM method for asymptotically stationary dependent sequences. Here we provide a bound for the excess risk  $V_\infty(f - \hat{g}_n) - \inf_{g \in \mathcal{G}} V_\infty(f - g)$ , where a control variate  $\hat{g}_n \in \mathcal{G}$  is chosen by minimization of the spectral variance  $V_n$  based on  $(X_k)_{k=0}^{n-1}$ , that is,  $\hat{g}_n \in \arg \min V_n(f - g)$ . The precise definition of  $V_n$  will be given in subsection 2.1. In section 3, we apply these results to Markov chains which are uniformly geometrically ergodic in Wasserstein distance. While subsection 3.1 is devoted to the (unadjusted) Langevin dynamics, in subsection 3.2 we use the ESVM approach for variance reduction in SGLD-type algorithms. We show that in both cases, the excess variance can be bounded, with high probability and up to logarithmic factors, as

$$V_\infty(f - \hat{g}_n) - \inf_{g \in \mathcal{G}} V_\infty(f - g) = O(n^{-1/2}).$$

This implies asymptotically valid confidence intervals (conditional on the sample used to learn  $\hat{g}_n$ ) of the form

$$\pi_N(f - \hat{g}_n) \pm \mathfrak{q} \sqrt{\frac{\inf_{g \in \mathcal{G}} V_\infty(f - g) + Cn^{-1/2}}{N}}$$

for some constant  $C > 0$ . Note that these intervals can be much tighter than ones in (1.2), provided that  $n$  is large and  $\inf_{g \in \mathcal{G}} V_\infty(f - g)$  is small. The latter condition is satisfied if the class  $\mathcal{G}$  is rich enough. In section 4, we illustrate performance of the proposed variance reduction method on various benchmark problems.

**Notations.** Let  $(\mathbf{X}, \mathbf{d})$  be a complete separable metric space. Define the Lipschitz norm of a real-valued function  $h$  by  $\|h\|_{\text{Lip}} := \sup_{x \neq y \in \mathbf{X}} \{|h(y) - h(x)| / \mathbf{d}(x, y)\}$ . We denote by  $\text{Lip}_{\mathbf{d}}(L)$  and  $\text{Lip}_{b, \mathbf{d}}(L, B)$  the class of Lipschitz (resp., bounded Lipschitz) functions on  $\mathbf{X}$  with  $\|h\|_{\text{Lip}} \leq L$  (resp.,  $\|h\|_{\text{Lip}} \leq L$  and  $|h|_\infty \leq B$ ). Further, let  $\mathbb{M}_1(\mathbf{X})$  be a set of probability measures on  $\mathbf{X}$ . We denote for  $p \geq 1$ ,  $\mathbb{S}_p(\mathbf{X}, \mathbf{d}) := \{\lambda \in \mathbb{M}_1(\mathbf{X}) : \int_{\mathbf{X}} \mathbf{d}^p(x, y) \lambda(dy) < \infty \text{ for all } x \in \mathbf{X}\}$ . For  $\lambda, \nu \in \mathbb{M}_1(\mathbf{X})$ , we denote their coupling set by  $\Pi(\lambda, \nu)$ , i.e.,  $\xi \in \Pi(\lambda, \nu)$  is the measure on  $\mathbf{X} \times \mathbf{X}$  satisfying for all  $A \in \mathcal{B}(\mathbf{X})$ ,  $\xi(A, \mathbf{X}) = \lambda(A)$  and  $\xi(\mathbf{X}, A) = \nu(A)$ . For  $p \geq 1$  and  $\lambda, \nu \in \mathbb{S}_p(\mathbf{X}, \mathbf{d})$ , let  $W_p^{\mathbf{d}}(\lambda, \nu) := \inf_{\Pi(\lambda, \nu)} \{\int_{\mathbf{X} \times \mathbf{X}} \mathbf{d}^p(x, y) \xi(dx, dy)\}^{1/p}$  be the Wasserstein distance of order  $p$  between  $\lambda$  and  $\nu$ . For  $\lambda, \nu \in \mathbb{M}_1(\mathbf{X})$ , let  $\text{KL}(\lambda|\nu)$  be the Kullback–Leibler divergence of  $\lambda$  with respect to  $\nu$ , i.e.,  $\text{KL}(\lambda|\nu) = \int \log(d\lambda/d\nu) d\lambda$  if  $\lambda \ll \nu$  and  $\text{KL}(\lambda|\nu) = \infty$  otherwise. Finally, unless otherwise specified, the symbol  $\lesssim$  stands for an inequality up to an absolute constant not depending on parameters of the problem.

## 2. Empirical spectral variance minimization (ESVM).

**2.1. Method.** Let  $(\Omega, \mathfrak{F}, (\mathfrak{F}_k)_{k \geq 0}, \mathbb{P})$  be a filtered probability space and  $(X_k)_{k=0}^\infty$  be a random process adapted to the filtration  $(\mathfrak{F}_k)_{k \geq 0}$  and taking values in  $\mathbf{X}$ . Let  $f : \mathbf{X} \rightarrow \mathbb{R}$  be a function such that  $\pi(f^2) < \infty$  and  $\mathbb{E}[f^2(X_k)] < \infty$  for all  $k \in \mathbb{N}$ . Let also  $\mathcal{G}$  be a set of control variates, that is, functions  $g \in \mathcal{G}$  satisfying  $\pi(g^2) < \infty$ ,  $\pi(g) = 0$ , and  $\mathbb{E}[g^2(X_k)] < \infty$  for all  $k \in \mathbb{N}$ . Particular examples of classes  $\mathcal{G}$  are given below in [section 3](#). Denote the class of functions  $h = f - g$  for  $g \in \mathcal{G}$  by  $\mathcal{H}$ ,

$$\mathcal{H} := \{f - g : g \in \mathcal{G}\}.$$

To shorten notation, we write  $\tilde{h} = h - \pi(h)$  for  $h \in \mathcal{H}$ .

We impose the following covariance stationarity condition on  $(X_k)_{k=0}^\infty$  to ensure that the asymptotic variance  $V_\infty(h)$  from [\(1.3\)](#) is well-defined for any  $h \in \mathcal{H}$ .

**(CS)** For any  $h \in \mathcal{H}$ , there exists a symmetric, summable, and positive semidefinite sequence  $(\rho^{(h)}(\ell))_{\ell \in \mathbb{Z}}$  satisfying the following conditions:

- (1)  $\rho^{(h)}(0) = \text{Var}_\pi(h)$ ;
- (2) for any  $\ell \in \mathbb{N}_0$  and constant  $R > 0$  independent of  $h$  and  $\ell$ ,
 
$$\sum_{k \in \mathbb{N}_0} |\mathbb{E}[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell)| \leq R;$$
- (3)  $\lim_{\ell \rightarrow \infty} \sum_{k \in \mathbb{N}_0} |\mathbb{E}[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell)| = 0$ .

**Proposition 2.1.** *Assume that the condition (CS) holds. Then, for all  $h \in \mathcal{H}$ , the asymptotic variance  $V_\infty(h)$  defined in [\(1.3\)](#) exists and can be represented as*

$$(2.1) \quad V_\infty(h) = \sum_{\ell \in \mathbb{Z}} \rho^{(h)}(\ell).$$

*Proof.* See [Appendix A.1](#). ■

The spectral variance estimator  $V_n(h)$  is based on truncation and weighting of the sample autocovariance functions:

$$(2.2) \quad V_n(h) := \sum_{|\ell| < b_n} w_n(\ell) \rho_n^{(h)}(\ell),$$

where  $w_n$  is the lag window,  $b_n$  is the truncation point, and  $\rho_n^{(h)}(\ell)$  is the sample autocovariance function given, for  $\ell \in \mathbb{N}_0$ , by

$$(2.3) \quad \rho_n^{(h)}(\ell) = \rho_n^{(h)}(-\ell) := n^{-1} \sum_{k=0}^{n-\ell-1} (h(X_k) - \pi_n(h))(h(X_{k+\ell}) - \pi_n(h)).$$

Here the truncation point  $b_n$  is an integer depending on  $n$ , and the lag window  $w_n$  is a kernel of the form  $w_n(\ell) = w(\ell/b_n)$ , where  $w$  is a symmetric nonnegative function supported on  $[-1, 1]$  such that  $\sup_{y \in [0,1]} |w(y)| \leq 1$  and  $w(y) = 1$  for  $y \in [-1/2, 1/2]$ . There are several other estimates for the asymptotic variance  $V_\infty(h)$ ; see [\[22\]](#) and the references therein. The ESVM estimator is obtained by

$$(2.4) \quad \hat{h}_n \in \arg \min_{h \in \mathcal{H}} V_n(h).$$

The ESVM method is summarized in [Algorithm 2.1](#).

**Algorithm 2.1.** ESVM method

- Input:** Two independent sequences:  $\mathbf{X}_n = (X_k)_{k=0}^{n-1}$  and  $\mathbf{X}'_N = (X'_k)_{k=0}^{N-1}$ .
1. Choose a class  $\mathcal{G}$  of functions with  $\pi(g) = 0$  for all functions  $g \in \mathcal{G}$ .
  2. Find  $\hat{g}_n \in \arg \min_{g \in \mathcal{G}} V_n(f - g)$ , where  $V_n$  is computed based on  $\mathbf{X}_n$ .
- Output:**  $\pi_N(f - \hat{g}_n)$  computed based on  $\mathbf{X}'_N$ .

**2.2. Theoretical analysis.** For our theoretical analysis, instead of looking for a function with the smallest spectral variance in the whole class  $\mathcal{H}$  we will perform optimization over a finite approximation (net) of  $\mathcal{H}$ . It turns out that both estimators have similar theoretical properties. Fix some  $\varepsilon > 0$ . Assuming that the class  $\mathcal{H}$  is totally bounded, let  $\mathcal{H}_\varepsilon$  be a minimal  $\varepsilon$ -net in the  $L^2(\pi)$ -norm, that is, the smallest possible (finite) collection of functions  $\mathcal{H}_\varepsilon \subset \mathcal{H}$  with the property that for any  $h \in \mathcal{H}$  there exists  $h_\varepsilon \in \mathcal{H}_\varepsilon$  such that the distance between  $h$  and  $h_\varepsilon$  in  $L^2(\pi)$ -norm is less than or equal to  $\varepsilon$ . The cardinality of  $\mathcal{H}_\varepsilon$  is called the covering number and is denoted by  $|\mathcal{H}_\varepsilon|$ . Define

$$\hat{h}_{n,\varepsilon} \in \arg \min_{h \in \mathcal{H}_\varepsilon} V_n(h).$$

To obtain a quantitative bound for the asymptotic variance of  $\hat{h}_{n,\varepsilon}$ , we need to specify the decay rate of the sequence  $(\rho^{(h)}(\ell))_{\ell \in \mathbb{Z}}$  from (CS).

(CD) There exist  $\varsigma > 0$  and  $\lambda \in [0, 1)$  such that, for any  $h \in \mathcal{H}$  and  $\ell \in \mathbb{N}_0$ ,

$$|\rho^{(h)}(\ell)| \leq \varsigma \lambda^\ell.$$

The following theorem provides a general bound on the excess of asymptotic variance.

**Theorem 2.2.** Assume that the conditions (CS) and (CD) hold. Assume additionally that for any  $n \in \mathbb{N}$  there exists a decreasing continuous function  $\alpha_n$  satisfying

$$\sup_{h \in \mathcal{H}} \mathbf{P}\left(|V_n(h) - \mathbf{E}[V_n(h)]| > t\right) \leq \alpha_n(t), \quad t > 0.$$

Then, for any  $\delta \in (0, 1)$  and  $\varepsilon > 0$ , it holds with probability at least  $1 - \delta$  that

$$\begin{aligned} V_\infty(\hat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty(h) &\lesssim \alpha_n^{-1}\left(\frac{\delta}{2|\mathcal{H}_\varepsilon|}\right) + (\sqrt{R}n^{-1/2} + \sqrt{D})b_n\varepsilon + \sqrt{RD}b_n n^{-1/2} \\ &\quad + (R + \varsigma(1 - \lambda)^{-1})b_n n^{-1} + \varsigma(1 - \lambda)^{-2}n^{-1} + \varsigma(1 - \lambda)^{-1}\lambda^{b_n/2}, \end{aligned}$$

where  $\alpha_n^{-1}$  is an inverse function for  $\alpha_n$  and  $D = \sup_{h \in \mathcal{H}} \text{Var}_\pi(h)$ .

*Proof.* See Appendix A.2. ■

Under some additional assumptions on the covering number of  $\mathcal{H}$  and the function  $\alpha_n(t)$ , a suitable choice of the size of  $\varepsilon$ -net and the truncation point  $b_n$  yields the following high-probability bound:

$$V_\infty(\hat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty(h) \lesssim n^{-1/(2+\rho)} \quad \text{for some } \rho > 0,$$

where  $\lesssim$  stands for inequality up to a constant depending on  $\lambda, R, D$ , and  $\varsigma$ . In the next section we shall apply Theorem 2.2 to the analysis of the ESVM algorithm for dependent sequences in the unadjusted Langevin algorithm (ULA) and SGLD.

**3. Applications.** In general, [Theorem 2.2](#) can be applied to different types of dependent sequences satisfying conditions [\(CS\)](#) and [\(CD\)](#). In what follows, we let  $(\mathbf{X}, \mathbf{d})$  be a complete separable metric space (equipped with its Borel  $\sigma$ -algebra  $\mathcal{X}$ ) and consider  $P$  to be a Markov kernel on  $(\mathbf{X}, \mathcal{X})$ . Let  $\Omega = \mathbf{X}^{\mathbb{N}}$  be the set of  $\mathbf{X}$ -valued sequences endowed with the  $\sigma$ -field  $\mathfrak{F} = \mathcal{X}^{\mathbb{N}}$ ,  $(X_k)_{k=0}^{\infty}$  be the coordinate process, and  $\mathfrak{F}_k = \sigma(X_\ell, \ell \leq k)$  be the canonical filtration. For every probability measure  $\xi$  on  $(\mathbf{X}, \mathcal{X})$  there exists a unique probability  $\mathbb{P}_\xi$  on  $(\mathbf{X}^{\mathbb{N}}, \mathcal{X}^{\otimes \mathbb{N}})$  such that the coordinate process  $(X_k)_{k=0}^{\infty}$  is a Markov chain with Markov kernel  $P$  and initial distribution  $\xi$ . We denote by  $\mathbb{E}_\xi$  the associated expectation. We focus below on the case where  $P$  is  $W_p^{\mathbf{d}}$ -uniformly ergodic for  $p = 1$  or  $p = 2$ .

**(WE)- $p$**  There exist  $x_0 \in \mathbf{X}$  such that  $\int_{\mathbf{X}} \mathbf{d}(x_0, x)P(x_0, dx) < \infty$  and a constant  $\Delta_p(P) \in [0, 1)$  such that

$$\sup_{(x, x') \in \mathbf{X}^2, x \neq x'} \frac{W_p^{\mathbf{d}}(\delta_x P, \delta_{x'} P)}{\mathbf{d}(x, x')} = \Delta_p(P).$$

[18, Theorem 20.3.4] shows that if **(WE)- $p$**  holds for some  $p \geq 1$ , then  $P$  admits a unique invariant probability measure which is denoted by  $\pi$  below. Moreover,  $\pi \in \mathbb{S}_p(\mathbf{X}, \mathbf{d})$  and for any  $\xi \in \mathbb{S}_p(\mathbf{X}, \mathbf{d})$ ,

$$(3.1) \quad W_p^{\mathbf{d}}(\xi P^n, \pi) \leq \Delta_p^n(P) W_p^{\mathbf{d}}(\xi, \pi), \quad n \in \mathbb{N}.$$

If there is no risk of confusion, we denote for simplicity  $\Delta_p = \Delta_p(P)$ . Let us start with a general result for Markov kernels satisfying **(WE)-2**. We show below that this assumption implies [\(CS\)](#) and [\(CD\)](#) when  $\mathcal{H}$  is a subset of Lipschitz functions, and establish an exponential concentration inequality for  $V_n(h)$ ,  $h \in \mathcal{H}$ . As was emphasized in [31] and [17], powerful tools for exploring concentration properties of  $W_2^{\mathbf{d}}$ -ergodic Markov kernels are the transportation cost-information inequalities.

**Definition 3.1.** For  $p \geq 1$ , we say that  $\mu \in \mathbb{M}_1(\mathbf{X})$  satisfies  $L^p$ -transportation cost-information inequality with constant  $\alpha > 0$  if for any  $\nu \in \mathbb{M}_1(\mathbf{X})$ ,  $W_p^{\mathbf{d}}(\mu, \nu) \leq \sqrt{2\alpha \text{KL}(\nu|\mu)}$ . We write briefly  $\mu \in T_p(\alpha)$  for this relation.

$L^p$ -transportation cost-information inequalities are well-studied in the literature; see, for instance, [4] and references therein. The cases  $p = 1$  and  $p = 2$  are of particular interest. Relations between  $T_1(\alpha)$  and concentration inequalities are covered in [29] and [8]. In particular,  $T_1(\alpha)$  is known to be equivalent to Gaussian concentration for all Lipschitz functions; see [8]. In turn  $T_2(\alpha)$  is a stronger inequality than  $T_1(\alpha)$ . It was first established for the standard Gaussian measure on  $\mathbb{R}^d$  by Talagrand in [44]. Moreover, the celebrated result by Bakry and Émery [3] implies that the measure  $\pi(dx) = e^{-U(x)} dx$  satisfies  $T_2(\alpha)$  if  $\nabla^2 U \geq \alpha^{-1} \mathbf{I}$ ; see [4, Chapter 9.6]. We are especially interested in  $T_2(\alpha)$ , since it is known to be stable under both independent and Markovian tensorizations; see [37] and [17].

Our results on  $W_2^{\mathbf{d}}$ -ergodic Markov kernels are summarized below.

**Proposition 3.2.** Let  $\mathcal{H} \subseteq \text{Lip}_{\mathbf{d}}(L)$ , and assume that **(WE)-2** holds. Then, for any initial distribution  $\xi \in \mathbb{S}_2(\mathbf{X}, \mathbf{d})$ , [\(CS\)](#) is satisfied with

$$(3.2) \quad \rho^{(h)}(\ell) = \mathbb{E}_\pi[\tilde{h}(X_0)\tilde{h}(X_{|\ell|})], \quad R = A_1 L^2 (1 - \Delta_2)^{-1} W_2(\xi, \pi),$$

where  $A_1$  is a constant given in (A.12), and (CD) is satisfied with

$$(3.3) \quad \varsigma = L\sqrt{D} \left[ \int \{W_2^d(\delta_x, \pi)\}^2 \pi(dx) \right]^{1/2}, \quad \lambda = \Delta_2, \quad D = \sup_{h \in \mathcal{H}} \text{Var}_\pi(h).$$

Moreover, if  $P(x, \cdot) \in T_2(\alpha)$  for any  $x \in \mathbf{X}$  and some  $\alpha > 0$ , then, for any initial distribution  $\xi \in T_2(\alpha)$ ,  $n \in \mathbb{N}$ , and  $t > 0$ ,

$$(3.4) \quad \mathbb{P}_\xi(|V_n(h) - \mathbb{E}_\xi[V_n(h)]| \geq t) \leq 2 \exp\left(-\frac{(1 - \Delta_2)^2 nt^2}{c\alpha L^2 b_n^2 (D + Rn^{-1} + t)}\right),$$

where  $c > 0$  is an absolute constant.

*Proof.* See Appendix A.3. ■

It is also possible to remove a quite restrictive assumption  $P(x, \cdot) \in T_2(\alpha)$  and to relax (WE)-2 to (WE)-1, but in this case (CS) and (CD) can be verified only for  $\mathcal{H}$  being a subset of bounded Lipschitz functions. As a price for such a generalization, the exponential concentration bound is replaced by a polynomial one.

**Proposition 3.3.** *Let  $\mathcal{H} \subset \text{Lip}_{b,d}(L, B)$ , and assume that (WE)-1 holds. Then for any initial distribution  $\xi \in \mathbb{S}_1(\mathbf{X}, d)$ , (CS) is satisfied with*

$$(3.5) \quad \rho^{(h)}(\ell) = \mathbb{E}_\pi[\tilde{h}(X_0)\tilde{h}(X_{|\ell|})], \quad R = A_2 B \left(1 - \Delta_1^{1/2}\right)^{-1},$$

where  $A_2$  is a constant given in (A.18), and (CD) is satisfied with

$$(3.6) \quad \varsigma = 2LB \int W_1^d(\delta_x, \pi)\pi(dx), \quad \lambda = \Delta_1, \quad D = \sup_{h \in \mathcal{H}} \text{Var}_\pi(h).$$

Moreover, for any  $p \in \mathbb{N}$ ,

$$(3.7) \quad \mathbb{P}_\xi(|V_n(h) - \mathbb{E}_\xi[V_n(h)]| \geq t) \leq \frac{C_{R,1}^p B^{2p} b_n^{3p/2} p^p}{n^{p/2} t^p} + \frac{C_{R,2}^p B^{2p} b_n^{2p} p^{2p}}{n^{p-1} t^p},$$

where constants  $C_{R,1}$  and  $C_{R,2}$  are given in (A.28).

*Proof.* See Appendix A.4. ■

**3.1. Langevin dynamics.** In this case,  $\mathbf{X} = \mathbb{R}^d$ , and we assume that  $\pi$  has an everywhere positive density w.r.t. the Lebesgue measure, i.e.,  $\pi(\theta) = Z^{-1}e^{-U(\theta)}$ , where  $Z = \int e^{-U(\vartheta)} d\vartheta$  is the normalization constant. Consider the first-order Euler–Maruyama discretization of the Langevin dynamics from (1.4),

$$(3.8) \quad \theta_{k+1} = \theta_k - \gamma \nabla U(\theta_k) + \sqrt{2\gamma} \xi_{k+1},$$

where  $\gamma > 0$  is a step size and  $(\xi_k)_{k=1}^\infty$  is an i.i.d. sequence of the standard Gaussian  $d$ -dimensional random vectors. The idea of using (3.8) to approximately sample from  $\pi$  has been advocated in [40], whose authors coined the term unadjusted Langevin algorithm (ULA). Consider the following assumption on  $U$ .

(ULA) The function  $U$  is continuously differentiable on  $\mathbb{R}^d$  with gradient  $\nabla U$  satisfying the following two conditions.

- (1) Lipschitz gradient: there exists  $L_U > 0$  such that for all  $\theta, \theta' \in \mathbb{R}^d$  it holds that  $\|\nabla U(\theta) - \nabla U(\theta')\| \leq L_U \|\theta - \theta'\|$ .
- (2) Strong convexity: there exists a constant  $m_U > 0$  such that for all  $\theta, \theta' \in \mathbb{R}^d$  it holds that  $U(\theta') \geq U(\theta) + \langle \nabla U(\theta), \theta' - \theta \rangle + (m_U/2)\|\theta' - \theta\|^2$ .

The ULA has been widely studied under the above assumptions; see, for example, [21] and [13]. As known from [21], under (ULA) the associated Markov kernel, denoted by  $P_\gamma^{(\text{ULA})}$ , is  $W_2^d$ -uniformly ergodic. For completeness, we state below the following proposition [21, Proposition 3].

**Proposition 3.4.** *Assume (ULA), and set  $\kappa = 2m_U L_U / (m_U + L_U)$ . Then for any step size  $\gamma \in (0, 2/(m_U + L_U))$ ,  $P_\gamma^{(\text{ULA})}$  satisfies (WE)-2 with  $d(\vartheta, \vartheta') = \|\vartheta - \vartheta'\|$  and  $\Delta_2 = \sqrt{1 - \kappa\gamma}$ . Moreover,  $P_\gamma^{(\text{ULA})}$  has a unique invariant measure  $\pi_\gamma^{(\text{ULA})}$ .*

It is shown in [21, Corollary 7] that, for any step size  $\gamma \in (0, 2/(m_U + L_U))$ ,

$$W_2^d(\pi, \pi_\gamma^{(\text{ULA})}) \leq \sqrt{2}\kappa^{-1/2} L_U \gamma^{1/2} \{\kappa^{-1} + \gamma\}^{1/2} \{2d + dL_U^2 \gamma / m_U + dL_U^2 \gamma^2 / 6\}^{1/2}.$$

We define the asymptotic variance as

$$V_\infty^{(\text{ULA})}(h) := \sum_{\ell \in \mathbb{Z}} \mathbb{E}_{\pi_\gamma^{(\text{ULA})}} \left[ (h(X_0) - \pi_\gamma^{(\text{ULA})}(f)) (h(X_{|\ell|}) - \pi_\gamma^{(\text{ULA})}(f)) \right].$$

At each iteration of the algorithm,  $\nabla U$  is computed. Hence it is an appealing option to use this gradient to construct Stein control variates (see, for instance, [1], [33], and [36]), given by

$$(3.9) \quad g_\phi(\theta) = -\langle \phi(\theta), \nabla U(\theta) \rangle + \text{div}(\phi(\theta)),$$

where  $\phi : \mathbf{X} \rightarrow \mathbb{R}^d$  is a continuously differentiable Lipschitz function,  $\langle \cdot, \cdot \rangle$  is the standard scalar product in  $\mathbb{R}^d$ , and  $\text{div}(\phi)$  is the divergence of  $\phi$ . Under rather mild conditions on  $\pi$  and  $\phi$ , it follows from the integration by parts that  $\pi(g_\phi) = 0$  (see [33, Propositions 1 and 2]). Note that if  $\phi(\theta) \equiv b$ ,  $b \in \mathbb{R}^d$ , we get  $g_b(\theta) = -\langle b, \nabla U(\theta) \rangle$ . Then for a parametric class  $\mathcal{H} = \{f - g_b : \|b\| \leq B\}$ , assuming that  $f \in \text{Lip}_d(L_1)$  and that condition (ULA) holds, we get  $\mathcal{H} \subset \text{Lip}_d(\max(L_1, BL_U))$ . For other approaches to construct control variates we refer the reader to [27], [16], and [9]. The next result follows now from Theorem 2.2 and Proposition 3.2.

**Theorem 3.5.** *Let  $\mathcal{H} \subset \text{Lip}_d(L)$ , and assume that (ULA) holds. Assume additionally that  $\xi \in T_2(\beta)$  for some  $\beta > 0$ . Fix any  $\gamma \in (0, 2/(m_U + L_U))$ , and set  $b_n = 2\lceil \log(n) / \log(1/\Delta_2) \rceil$  with  $\Delta_2 = \sqrt{1 - \kappa\gamma}$  and  $\kappa = 2m_U L_U / (m_U + L_U)$ . Then, for any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*



$$V_\infty^{(\text{ULA})}(\widehat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty^{(\text{ULA})}(h) \lesssim C_1 \varepsilon \log(n) + C_2 \sqrt{\frac{\log^2(n) \log(|\mathcal{H}_\varepsilon|/\delta)}{n}} + C_3 \frac{\log^2(n) \log(|\mathcal{H}_\varepsilon|/\delta)}{n},$$

where

$$C_1 = \frac{\sqrt{R} + \sqrt{D}}{\kappa\gamma}, C_2 = \frac{L\sqrt{(\beta \vee \gamma)(D + R)}}{\kappa^2\gamma^2} + \frac{\sqrt{DR}}{\kappa\gamma}, C_3 = \frac{L^2(\beta \vee \gamma)}{\kappa^4\gamma^4} + \frac{R}{\kappa\gamma} + \frac{\varsigma}{\kappa^2\gamma^2}$$

with  $R, \varsigma$  from Proposition 3.2 and  $D = \sup_{h \in \mathcal{H}} \text{Var}_{\pi_\gamma^{(\text{ULA})}}(h)$ .

*Proof.* The Markov kernel associated to ULA can be written as  $P_\gamma^{(\text{ULA})}(\theta, \cdot) = \mathcal{N}(\theta - \gamma \nabla U(\theta), 2\gamma I_d)$ . Hence, by [4, Theorem 9.2.1],  $P_\gamma^{(\text{ULA})}(\theta, \cdot) \in \mathbb{T}_2(2\gamma)$  for any  $\gamma > 0$ . By Proposition 3.4, (WE) holds with  $\Delta_2 = \sqrt{1 - \gamma\kappa}$ . Hence Proposition 3.2 applies with  $\alpha = 2(\beta \vee \gamma)$ . Direct computation of the inverse function in the right-hand side of (3.4) leads to

$$\alpha_n^{-1} \left( \frac{\delta}{2|\mathcal{H}_\varepsilon|} \right) \leq \frac{4b_n^2 L^2(\beta \vee \gamma) \log(4|\mathcal{H}_\varepsilon|/\delta)}{(1 - \Delta_2)^2 n} + \frac{4b_n L \sqrt{(\beta \vee \gamma)(D + R)} \log(4|\mathcal{H}_\varepsilon|/\delta)}{(1 - \Delta_2) \sqrt{n}}. \quad \blacksquare$$

**Corollary 3.6.** *Under the assumptions of Theorem 3.5, the following holds.*

- (1) *If class  $\mathcal{H}$  is parametric, that is,  $|\mathcal{H}_\varepsilon| \leq C_\rho \varepsilon^{-\rho}$  for all  $\varepsilon \in (0, 1)$  and some constants  $C_\rho, \rho > 0$ , then it holds for any  $\varepsilon \in (0, 1/\sqrt{n})$  with probability at least  $1 - 1/n$ ,*

$$V_\infty^{(\text{ULA})}(\widehat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty^{(\text{ULA})}(h) \lesssim n^{-1/2} \log^{1/2}(n).$$

- (2) *If class  $\mathcal{H}$  is nonparametric, that is,  $|\mathcal{H}_\varepsilon| \leq C_\rho \exp(\varepsilon^{-\rho})$  for all  $\varepsilon \in (0, 1)$  and some constants  $C_\rho, \rho > 0$ , then it holds for any  $\varepsilon \in (0, 1/\sqrt{n})$  with probability at least  $1 - 1/n$ ,*

$$V_\infty^{(\text{ULA})}(\widehat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty^{(\text{ULA})}(h) \lesssim n^{-1/(2+\rho)}.$$

Here  $\lesssim$  stands for inequality up to a constant depending on  $\rho$  and other constants from Theorem 3.5. Moreover, if additionally the constant  $\pi_\gamma^{(\text{ULA})}(f)$  is in the class  $\mathcal{H}$ , then  $\inf_{h \in \mathcal{H}} V_\infty^{(\text{ULA})}(h) = 0$  and these bounds hold for the asymptotic variance itself.

*Discussion.* It is well-known that if  $\hat{f}$  satisfies the so-called Poisson equation  $P_\gamma^{(\text{ULA})} \hat{f} - \hat{f} = -f + \pi_\gamma^{(\text{ULA})}(f)$ , then by taking  $g^* = \hat{f} - P_\gamma^{(\text{ULA})} \hat{f}$  as a control variate, we get  $\pi_\gamma^{(\text{ULA})}(f - g^*) = \pi_\gamma^{(\text{ULA})}(f)$  and  $V_\infty^{(\text{ULA})}(f - g^*) = 0$ . The property  $h^* = f - g^* = \pi_\gamma^{(\text{ULA})}(f) \in \mathcal{H}$  can be achieved by taking, for example,  $\mathcal{H}$  to be a ball in a Sobolev space. Namely, let  $W_2^s = \{h \in L^2(\lambda) : D^\alpha h \in L^2(\lambda), \forall |\alpha| \leq s\}$  be the Sobolev space; here  $\lambda$  is the Lebesgue measure on  $\mathbb{R}^d$ ,  $\alpha = (\alpha_1, \dots, \alpha_d)$  is a multi-index with  $|\alpha| = \alpha_1 + \dots + \alpha_d$ , and  $D^\alpha$  stands for the differential operator  $D^\alpha = \partial^{|\alpha|} / \partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}$ . The weighted Sobolev space  $W_2^s(\langle x \rangle^\beta)$ ,  $\beta \in \mathbb{R}$ , for a polynomial weighting function  $\langle x \rangle^\beta = (1 + \|x\|^2)^{\beta/2}$  is defined by  $W_2^s(\langle x \rangle^\beta) = \{h : h \cdot \langle x \rangle^\beta \in W_2^s\}$ . Let  $\mathcal{H}$  be a norm-bounded subset of  $W_2^s(\langle x \rangle^\beta)$  with  $\beta \in \mathbb{R}$  and  $s - d/2 > 0$ . Suppose also that  $\|\langle x \rangle^{\alpha-\beta}\|_{L^2(\pi_\gamma^{(\text{ULA})})} < \infty$  for some  $\alpha > 0$ . Then  $|\mathcal{H}_\varepsilon| \lesssim \exp(\varepsilon^{-d/s})$ , provided

that  $\alpha > s - d/2$ ; see [35, Corollary 4]. Note that  $h^* = \pi_\gamma^{(\text{ULA})}(f) \in W_2^s(\langle x \rangle^\beta)$  for any  $s > 0$  and any  $\beta < -1$  so that we can take  $\mathcal{H}$  as a norm-bounded subset of  $W_2^s(\langle x \rangle^\beta)$  for arbitrary large  $s > 0$ . Since  $\pi_\gamma^{(\text{ULA})}$  and all its derivatives have exponentially decaying tails (see [32]),  $\|\langle x \rangle^{\alpha-\beta}\|_{L^2(\pi_\gamma^{(\text{ULA})})} < \infty$  for any  $\alpha > 0$ , and one can achieve that  $|\mathcal{H}_\varepsilon| \lesssim \exp(\varepsilon^{-\delta})$  for arbitrary small  $\delta > 0$  and at the same time  $h^* \in W_2^s(\langle x \rangle^\beta)$ . Practically one can use Stein control variates of the form (3.9) with infinitely smooth and compactly supported functions  $\phi$ . This will guarantee that  $f - g_\phi \in W_2^s(\langle x \rangle^\beta)$  for some  $s > 0$ , provided that  $U$  is smooth enough and  $f \in W_2^s(\langle x \rangle^\beta)$ .

**3.2. Extension to the SGLD.** In this section, we shall consider the situations where the target  $\pi$  is given by the posterior distribution in the Bayesian inference problem, that is,  $\pi(\theta) \propto \exp(-U(\theta))$ , where  $U(\theta) = U_0(\theta) + \sum_{i=1}^K U_i(\theta)$  with  $K$  being a number of observations. Computing  $\nabla U(\theta)$  requires a computational budget that scales linearly with  $K$ . Hence it is often impossible to apply procedures based on discretization of Langevin dynamics directly. One possible solution advocated by [46] is to replace  $\nabla U(\theta)$  by an unbiased estimate. This gives rise to the SGLD algorithm, where the parameters are updated according to

$$(3.10) \quad \begin{aligned} \theta_{k+1} &= \theta_k - \gamma G(\theta_k, S_{k+1}) + \sqrt{2\gamma} \xi_{k+1}, \\ G(\theta, S) &= \nabla U_0(\theta) + KM^{-1} \sum_{i \in S} \nabla U_i(\theta), \end{aligned}$$

where each  $S_{k+1}$  is a random batch taking values in  $\mathfrak{S}_M$  (here  $\mathfrak{S}_M$  is the set of all subsets  $S$  of  $\{1, \dots, K\}$  with  $|S| = M$ ) which is sampled from a uniform distribution over  $\mathfrak{S}_M$  independently of  $\mathcal{F}_k$  (here  $(\mathcal{F}_k)_{k \geq 0}$  is the filtration generated by  $\{(\theta_\ell, S_\ell)\}_{\ell \geq 0}$ ). Note that  $\mathbb{E}[G(\theta_k, S_{k+1}) | \mathcal{F}_k] = \nabla U(\theta_k)$  and therefore  $G(\theta_k, S_{k+1})$  is an unbiased estimate of  $\nabla U(\theta_k)$ . The available variance reduction techniques for SGLD usually replace the stochastic gradient in (3.10) with more sophisticated estimates which preserve unbiasedness but have lower variance.

The simplest variance reduction technique is the fixed-point method (SGLD-FP) proposed in [2]. This method is applicable when the posterior distribution is strongly log-concave. We set  $\hat{\theta} \in \Theta$  to be a fixed value of the parameter, typically chosen to be close to the mode of posterior distribution. We estimate the gradient  $\nabla U(\theta)$  by

$$(3.11) \quad G_{\text{FP}}(\theta, S) = \nabla U_0(\theta) + KM^{-1} \sum_{i \in S} (\nabla U_i(\theta) - \nabla U_i(\hat{\theta})) + \sum_{i=1}^K \nabla U_i(\hat{\theta}).$$

The SGLD-FP algorithm is obtained by plugging this approximation into (3.10).

More sophisticated variance reduction methods typically use reference values  $(g_k^i)_{i=1}^K$  of the gradient  $(\nabla U_i)_{i=1}^K$  from previous iterates (and not only the last iterate); as a result, constructed sequence  $(\theta_k)_{k=0}^\infty$  is often not Markovian. One particular example is SAGALD method, adapted from [41], [15]. If  $i \in S_k$ , the reference value is updated, that is,  $g_{k+1}^i = \nabla U_i(\theta_k)$ . Otherwise, the reference value is simply propagated, that is,  $g_{k+1}^i = g_k^i$ . One then considers the following gradient estimator:

$$(3.12) \quad G_{\text{SAGA}}^k(\theta, S) = \nabla U_0(\theta) + KM^{-1} \sum_{i \in S} (\nabla U_i(\theta) - g_k^i) + g_k, \quad g_k = \sum_{i=1}^K g_k^i.$$

The recursion is initialized with  $g_0^i = \nabla U_i(\theta_0)$ ,  $i \in \{1, \dots, K\}$ , and  $g_0 = \sum_{i=1}^K g_0^i$ . Finally, the gradient is computed according to (3.12) and plugged into (3.10).

For theoretical analysis of SGLD and SGLD-FP algorithms we need the following assumptions on  $U$ . Without loss of generality, we consider only SGLD; the same reasoning applies to SGLD-FP.

**(SGLD)** The function  $U(\theta) = U_0(\theta) + \sum_{i=1}^K U_i(\theta)$  satisfies the following conditions.

- (1) Lipschitz gradient: for any  $i \in \{0, \dots, K\}$ ,  $U_i$  is continuously differentiable on  $\mathbb{R}^d$  with  $\tilde{L}_U$ -Lipschitz gradient.
- (2) Convexity: for any  $i \in \{0, \dots, K\}$ ,  $U_i$  is convex.
- (3) Strong convexity: there exists a constant  $m_U > 0$  such that for any  $\theta, \theta' \in \mathbb{R}^d$  it holds that  $U(\theta') \geq U(\theta) + \langle \nabla U(\theta), \theta' - \theta \rangle + (m_U/2)\|\theta' - \theta\|^2$ .

Note that using Stein control variates with SGLD-based sampling procedure (3.10) eliminates benefits of using  $G(\theta, S)$  instead of exact gradient  $\nabla U(\theta)$ . Following [23], we replace  $\nabla U$  by its stochastic counterpart. More precisely, for the  $k$ th iteration of SGLD algorithm, we consider the control variates of the form

$$(3.13) \quad g_\phi(\theta, S) = -\langle \phi(\theta), G(\theta, S) \rangle + \text{div}(\phi(\theta)).$$

The control variate  $g_\phi$  depends now on the pair  $(\theta, S)$ . Let  $\mathcal{H} = \{f(\theta) - g_\phi(x) : \phi \in \Phi\}$ , where  $x = (\theta, S) \in \mathbf{X} = \Theta \times \mathbf{S}_M$ . Consider another sequence  $(\tilde{S}_k)_{k=0}^\infty$  of independent batches uniformly distributed over  $\mathbf{S}_M$  such that for any  $k$ ,  $\tilde{S}_k$  is independent of  $\mathcal{F}_k$ . Denote by  $P_{\text{SGLD}}$  the transition kernel of SGLD, and let  $\Upsilon_M$  be a uniform distribution over  $\mathbf{S}_M$ . Set  $\bar{P} := P_{\text{SGLD}} \otimes \Upsilon_M$  and  $X_k = (\theta_k, \tilde{S}_k)$ .

**Proposition 3.7.** Assume (SGLD). Then for any step size  $\gamma \in (0, \tilde{L}_U^{-1}(K + 1)^{-1})$ ,  $\bar{P}$  satisfies (WE)-2 with  $\Delta_2 = \sqrt{1 - \gamma m_U}$  and  $d(x, x') = \|\vartheta - \vartheta'\| + \mathbb{1}_{\{S \neq S'\}}$  for any  $x = (\vartheta, S)$  and  $x' = (\vartheta', S')$ . Moreover,  $\bar{P}$  has a unique invariant measure  $\bar{\pi} = \pi_\gamma^{(\text{SGLD})} \otimes \Upsilon_M$ .

*Proof.* See Appendix A.5. ■

Similarly to Langevin dynamics, we define

$$V_\infty^{(\text{SGLD})}(h) := \sum_{\ell \in \mathbb{Z}} \mathbb{E}_{\bar{\pi}} \left[ (h(X_0) - \bar{\pi}(f)) (h(X_{|\ell|}) - \bar{\pi}(f)) \right].$$

**Theorem 3.8.** Let  $\mathcal{H} \subseteq \text{Lip}_{b,d}(L, B)$ , and assume that (SGLD) holds. Fix any  $\gamma \in (0, \tilde{L}_U^{-1}(K + 1)^{-1})$ , and set  $b_n = 2 \lceil \log(n) / \log(1/\Delta_1) \rceil$  with  $\Delta_1 = \sqrt{1 - \gamma m_U}$ . Then, for any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} & V_\infty^{(\text{SGLD})}(\hat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty^{(\text{SGLD})}(h) \\ & \lesssim C_4 \varepsilon \log(n) + C_5 \sqrt{\frac{\log^5(n)}{n}} \left( \frac{|\mathcal{H}_\varepsilon|}{\delta} \right)^{1/\log(n)} + C_6 \frac{\log n}{n}, \end{aligned}$$

where

$$C_4 = \frac{\sqrt{R} + \sqrt{D}}{m_U \gamma}, \quad C_5 = \frac{B^2 R_1(L, \xi)}{(m_U \gamma)^2} + \frac{B^2 R_2(L, \xi)}{(m_U \gamma)^{4+2/\log n}} + \frac{\sqrt{RD}}{m_U \gamma}, \quad C_6 = \frac{D(m_U \gamma) + \varsigma}{(m_U \gamma)^2}$$

with  $R, \varsigma$  from [Proposition 3.3](#),  $D = \sup_{h \in \mathcal{H}} \text{Var}_{\pi_\gamma^{(\text{SGLD})}}(h)$ , and constants  $R_1(L, \xi), R_2(L, \xi)$  which can be tracked from [\(A.27\)](#).

*Proof.* By [Proposition 3.7](#), (WE)-2 holds with  $\Delta_2 = \sqrt{1 - \gamma m_U}$ , and, by Lyapunov inequality, (WE)-1 also holds with  $\Delta_1 = \Delta_2$ . Hence, the second part of [Proposition 3.3](#) can be applied with  $p = \log n$ . The remaining part follows from [Theorem 2.2](#) with computation of the inverse function in the right-hand side of [\(3.7\)](#). ■

**Corollary 3.9.** *Under the assumptions of [Theorem 3.8](#), if class  $\mathcal{H}$  is parametric, that is,  $|\mathcal{H}_\varepsilon| \leq C_\rho \varepsilon^{-\rho}$  for all  $\varepsilon \in (0, 1)$  and some constants  $C_\rho, \rho > 0$ . Then for any  $\varepsilon \in (0, 1/\sqrt{n})$  it holds with probability at least  $1 - 1/n$  that*

$$V_\infty^{(\text{SGLD})}(\hat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty^{(\text{SGLD})}(h) \lesssim n^{-1/2} \log^{5/2}(n),$$

where  $\lesssim$  stands for inequality up to a constant depending on  $\rho$  and other constants from [Theorem 3.8](#). Moreover, if additionally  $\bar{\pi}(f) \in \mathcal{H}$ , then  $\inf_{h \in \mathcal{H}} V_\infty^{(\text{SGLD})}(h) = 0$ , and these bounds hold for the asymptotic variance itself.

**Remark 3.10.** If the class  $\mathcal{H}$  is constructed using Stein control variates, we can ensure the inclusion  $\mathcal{H} \subseteq \text{Lip}_{b,d}(L, B)$  by taking smooth and compactly supported functions  $\phi$ . This in turn can be achieved by multiplying a given smooth function  $\phi$  with a mollifier function, that is, an infinitely smooth compactly supported function.

**4. Experiments.** In this section, we numerically compare the following two methods to choose control variates: the EVM method, where a control variate is determined by minimizing the marginal variance (see [\(1.1\)](#)), and the ESVM method, where a control variate is determined by minimizing the spectral variance (see [\(2.2\)](#)). Implementation is available at [https://github.com/svsamsonov/vr\\_sg\\_mcmc](https://github.com/svsamsonov/vr_sg_mcmc).

**4.1. Toy example.** We first consider a multimodal distribution in  $\mathbb{R}^2$  from [\[38\]](#). Namely, let  $\pi(x_1, x_2) = Z^{-1} e^{-U(x_1, x_2)}$ , where  $Z$  is the normalization constant and

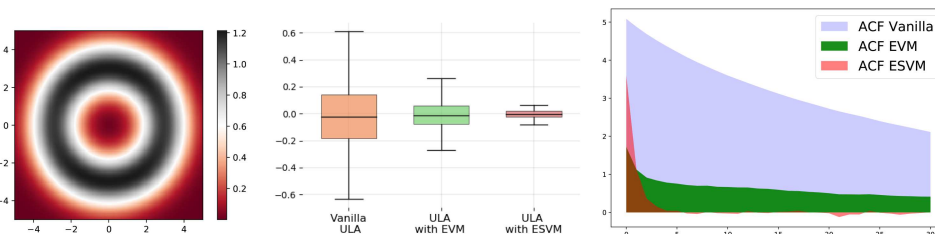
$$U(x_1, x_2) = \frac{(\|x\| - \mu)^2}{2M^2} - \log\left(e^{-(x_1 - \mu)^2/2\sigma^2} + e^{-(x_1 + \mu)^2/2\sigma^2}\right).$$

We choose  $M = 1$  and  $\mu = \sigma = 3$ ; the respective density profile is presented in [Figure 1](#). Our aim is to estimate  $\pi(f)$  with  $f(x_1, x_2) = x_1 + x_2$  using ULA. The parametric class  $g_\varphi$  in [\(3.9\)](#) is generated by  $\varphi(x) = \sum_{k=1}^p \beta_k \psi_k(x)$ , where  $\psi_k = e^{-\|x - \mu_k\|^2/2\sigma_\psi^2}$  with all  $\mu_k$  regularly spaced in  $[-3, 3] \times [-3, 3]$  and  $\sigma_\psi = 2$ . Details on the step size  $\gamma$  of the ULA, length of the burn-in period and test trajectories are summarized in [Table 1](#). Boxplots displaying variation of 100 estimates for EVM and ESVM are presented in the same [Figure 1](#). Furthermore, we compute sample autocovariance functions for a trajectory with and without adding ESVM and EVM control variates. The results reflect a spectacular decrease in high-order autocovariance for ESVM; see [Figure 1](#). Note that EVM aims at minimizing only the lag-zero autocovariance; that is why the autocovariance function for ESVM-adjusted trajectory decreases much faster.

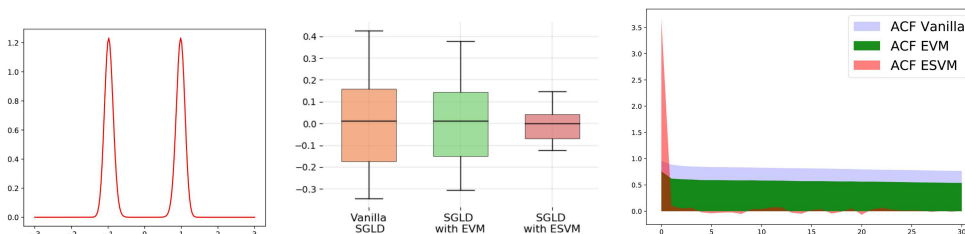
**4.2. Gaussian mixture model.** We consider posterior mean estimation for unknown parameter  $\mu$  in a Bayesian setup with normal prior  $\mu \sim \mathcal{N}(0, \sigma_\mu^2)$ ,  $\sigma_\mu^2 = 100$ , and sample  $(X_k)_{k=0}^{K-1}$ ,  $K = 100$ , drawn from the Gaussian mixture model

**Table 1**  
Experimental hyperparameters.

Experiment	$n_{\text{burn}}$	$n_{\text{test}}$	$\gamma$	Batch size
Toy example, subsection 4.1	$10^3$	$10^4$	0.1	-
Gaussian Mixture, subsection 4.2	$10^4$	$10^5$	0.01	10



**Figure 1.** Toy example from subsection 4.1. From left to right: (1) density profile, (2) boxplots displaying variation of 100 estimates for vanilla ULA, ULA with EVM, and ULA with ESVM, (3) sample autocovariance functions (ACFs) for a trajectory with and without ESVM and EVM.



**Figure 2.** Gaussian mixture model from subsection 4.2. From left to right: (1) density of the posterior distribution, (2) boxplots displaying variation of 100 estimates for vanilla SGLD, SGLD with EVM, and SGLD with ESVM, (3) sample autocovariance functions (ACFs) for a trajectory with and without ESVM and EVM.

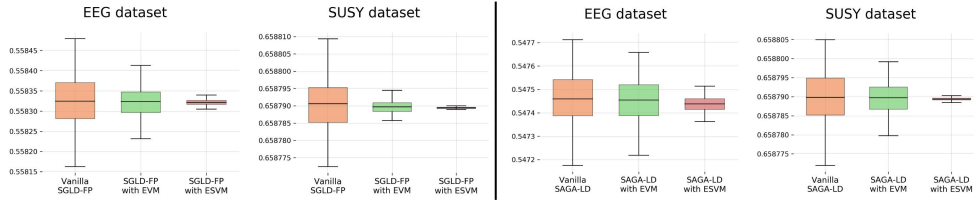
$$0.5\mathcal{N}(-\mu, \sigma^2) + 0.5\mathcal{N}(\mu, \sigma^2) \quad \text{with } \mu = 1, \sigma^2 = 1.$$

The density of the posterior distribution over  $\mu$  is given in Figure 2. It has 2 modes roughly corresponding to  $\mu = 1$  and  $\mu = -1$ . To generate data from this posterior distribution and estimate posterior mean, we use SGLD. The parametric class  $g_\varphi$  in (3.13) is generated by  $\varphi(x) = \beta_0 x^2 + \beta_1 x + \beta_2$ . Boxplots displaying variation of 100 estimates for EVM and ESVM and respective sample autocovariance functions are also presented in Figure 2. Note that the increase in lag-zero autocovariance for ESVM is explained by the additional randomness in (3.13). On contrary, EVM favors far too small coefficients to overcome this additional randomness, which leads to poor variance reduction.

**4.3. Bayesian logistic regression.** The probability of the  $i$ th output  $y_i \in \{-1, 1\}$ ,  $i = 1, \dots, K$ , is given by  $p(y_i | \mathbf{x}_i, \theta) = (1 + e^{-y_i \langle \theta, \mathbf{x}_i \rangle})^{-1}$ , where  $\mathbf{x}_i$  is a  $d \times 1$  vector of predictors and  $\theta$  is the vector of unknown regression coefficients. We complete the Bayesian model by considering the Zellner  $g$ -prior  $\mathcal{N}_d(0, g(\mathbf{X}^\top \mathbf{X})^{-1})$  for  $\theta$  where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$  is an  $K \times d$  design matrix; see [26, section 2]. Normalizing the covariates, for  $\tilde{\mathbf{x}}_i = (\mathbf{X}^\top \mathbf{X})^{-1/2} \mathbf{x}_i$  and  $\tilde{\theta} = (\mathbf{X}^\top \mathbf{X})^{1/2} \theta$ , we get  $\langle \theta, \mathbf{x}_i \rangle = \langle \tilde{\theta}, \tilde{\mathbf{x}}_i \rangle$ , under the Zellner  $g$ -prior,  $\tilde{\theta} \sim \mathcal{N}_d(0, g\mathbf{I}_d)$ . In our

**Table 2**  
*Experimental hyperparameters.*

Experiment	$n_{\text{burn}}$	$n_{\text{train}}$	$n_{\text{test}}$	$\gamma$	Batch size
Logistic regression, EEG dataset	$10^4$	$10^4$	$10^5$	0.1	15
Logistic regression, SUSY dataset	$10^5$	$10^5$	$10^6$	0.1	50



**Figure 3.** Bayesian logistic regression for EEG and SUSY datasets from subsection 4.3. Boxplots displaying variation of 100 estimates of average predictive distribution for (1) left panel: vanilla SGLD-FP, SGLD-FP with EVM, and SGLD-FP with ESVM, (2) right panel: vanilla SAGA-LD, SAGA-LD with EVM, and SAGA-LD with ESVM.

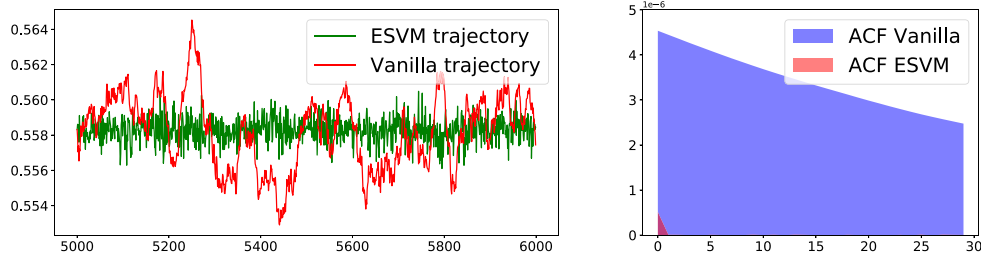
experiments we have not noticed significant impact of particular  $g$  value on EVM and ESVM performance and used  $g = 100$  as a default choice.

We analyze the performance of EVM and ESVM methods on two datasets from the UCI repository. The first dataset, EEG,<sup>1</sup> contains  $K = 14980$  observations in dimension  $d = 15$ ; the second dataset, SUSY,<sup>2</sup> has  $K = 500000$  observations in dimension  $d = 19$ . The data is first split into a training set  $\mathcal{T}_N^{\text{train}} = \{(y_i, \mathbf{x}_i)\}_{i=1}^K$  and a test set  $\mathcal{T}_K^{\text{test}} = \{(y'_i, \mathbf{x}'_i)\}_{i=1}^K$  by randomly picking  $K = 100$  test points from the data. We use the SGLD-FP and SAGA-LD algorithms to approximately sample from the posterior distribution  $p(\tilde{\theta} | \mathcal{T}_N^{\text{train}})$ . Given a sample  $(\tilde{\theta}_k)_{k=0}^{n-1}$ , we can estimate the predictive distribution for a fixed test point  $(y', \mathbf{x}')$ , that is,  $p(y' | \mathbf{x}') = \int_{\mathbb{R}^d} p(y' | \mathbf{x}', \tilde{\theta}) p(\tilde{\theta} | \mathcal{T}_N^{\text{train}}) d\tilde{\theta}$ , by computing the ergodic mean  $n^{-1} \sum_{k=0}^{n-1} f(\tilde{\theta}_k)$  for  $f(\tilde{\theta}) = p(y' | \mathbf{x}', \tilde{\theta})$ . To get rid of randomness caused by the random choice of a test point, we estimate the average predictive distribution for the whole test set  $\mathcal{T}_K^{\text{test}}$  by computing the ergodic mean for the function  $f(\tilde{\theta}) = K^{-1} \sum_{i=1}^K p(y'_i | \mathbf{x}'_i, \tilde{\theta})$ . Details on the step size  $\gamma$ , length of the burn-in period and test trajectories, and batch size are summarized in Table 2. Boxplots for the estimation of average predictive distribution are shown in Figure 3. Note that ESVM leads to a significant variance reduction for both SGLD-FP and SAGA-LD.

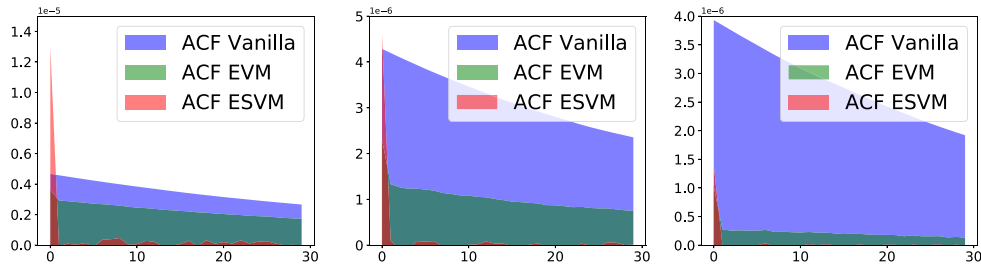
Further, for the EEG dataset we plot in Figure 4 a part of the trajectory  $f(\tilde{\theta}_m) = K^{-1} \sum_{i=1}^K p(y'_i | \mathbf{x}'_i, \tilde{\theta}_m)$  for 500 consecutive sample values  $\tilde{\theta}_m$  with and without adding the ESVM control variate. These trajectories are accompanied by the sample autocovariance functions for vanilla and variance-reduced samples for both EVM and ESVM. Again, since EVM aims at minimizing only lag-zero autocovariance, the decrease in autocovariance function for this method is smaller than for ESVM. We also report in Figure 5 how autocovariance functions change with batch sizes. Note that for small batch sizes ESVM still manages to remove correlations, while EVM almost fails.

<sup>1</sup><https://archive.ics.uci.edu/ml/datasets/EEG+Eye+State>.

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/SUSY>.



**Figure 4.** Bayesian logistic regression for the EEG dataset from subsection 4.3 with the batch size 150. From left to right: (1) part of a trajectory with and without ESVM, (2) sample autocovariance functions (ACFs) for a trajectory with and without ESVM.



**Figure 5.** Bayesian logistic regression for the EEG dataset from subsection 4.3. Comparison of sample autocovariance functions (ACFs) for different batch sizes. From left to right: batch size 5, 10, 50, respectively.

**4.4. Bayesian probabilistic matrix factorization.** A typical problem in recommendation systems is to predict a user’s rating for a particular item given other users’ ratings of this item and how a given user evaluated other items. A common approach to this problem is probabilistic matrix factorization via Bayesian inference; see [43]. Namely, we are interested in approximating matrix  $R \in \mathbb{R}^{M \times N}$ , where  $M$  is a number of users,  $N$  is a number of rated items, and  $R_{i,j}$  stands for the rating assigned by the  $i$ th user to the  $j$ th item. Due to natural limitations (user is unlikely to rate all possible items), we observe only a some small subset of elements of  $R$  and want to predict ratings of the hidden part. In probabilistic matrix factorization, we aim at representing  $R$  as a product  $R = U^T V + C$ , where  $U \in \mathbb{R}^{D \times M}$ ,  $V \in \mathbb{R}^{D \times N}$ , and  $C \in \mathbb{R}^{M \times N}$  is a matrix of biases with elements  $C_{i,j} = a_i + b_j$ ,  $a \in \mathbb{R}^M$ ,  $b \in \mathbb{R}^N$ . In the subsequent experiments we assume that rank parameter  $D = 10$  is fixed. The naive solution would be to find

$$U, V, a, b = \arg \min_{U, V, a, b} \sum_{(i,j) \in I_{\text{train}}} (R_{i,j} - \langle U_i, V_j \rangle - a_i - b_j)^2,$$

where  $I_{\text{train}}$  is a train subset of ratings. Unfortunately, optimizing this criterion leads to a significantly overfitted model. One possible approach to overcome overfitting is to consider the penalized model

$$U, V, a, b = \arg \min_{U, V, a, b} \sum_{(i,j) \in I_{\text{train}}} (R_{i,j} - \langle U_i, V_j \rangle - a_i - b_j)^2 + \lambda_U \|U\|^2 + \lambda_V \|V\|^2 + \lambda_a \|a\|^2 + \lambda_b \|b\|^2,$$

but it requires careful tuning of penalization coefficients  $\lambda_U, \lambda_V, \lambda_a, \lambda_b$ . We thus would benefit a lot from a Bayesian approach for tuning weights; this was pointed out in [43]. We follow a slightly simplified formulation proposed by [12]; that is, we consider

$$\begin{aligned} \lambda_U, \lambda_V, \lambda_a, \lambda_b &\sim \Gamma(1, 1), \quad U_{k,i} \sim \mathcal{N}(0, \lambda_U^{-1}), \quad V_{k,j} \sim \mathcal{N}(0, \lambda_V^{-1}), \\ a_i &\sim \mathcal{N}(0, \lambda_a^{-1}), \quad b_i \sim \mathcal{N}(0, \lambda_b^{-1}), \quad R_{i,j}|U, V \sim \mathcal{N}(\langle U_i, V_j \rangle + a_i + b_j, \tau^{-1}). \end{aligned}$$

In order to sample from the posterior distribution which we denote by  $p(\Theta|R)$ , where  $\Theta = \{U, V, a, b, \lambda_U, \lambda_V, \lambda_a, \lambda_b\}$ , we use the following two-step procedure:

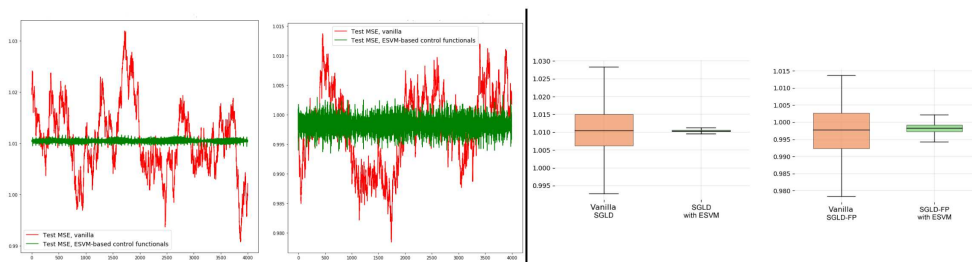
1. Sample from  $p(U, V, a, b|R, \lambda_U, \lambda_V, \lambda_a, \lambda_b)$  using SGLD or SGLD-FP with a minibatch size of 5000 observations with a step size  $\gamma = 10^{-4}$ . Sample for 1000 steps before updating the weights  $\lambda_U, \lambda_V, \lambda_a, \lambda_b$ ;
2. Sample new  $\lambda$  from  $p(\lambda_U, \lambda_V, \lambda_a, \lambda_b|U, V, a, b)$  using the Gibbs sampler.

The experiments are performed on the MovieLens dataset *ml-100k* ([link to dataset](#)). We apply our control variates procedure as a postprocessing step following [2]. The functional of interest is the mean squared error over the test subsample,  $f(U, V, a, b) = \sum_{(i,j) \in I_{\text{test}}} (R_{i,j} - \langle U_i, V_j \rangle - a_i - b_j)^2$ . Since the dimension of parameter space is very high, first-order control variates are the only option among Stein's control variates. Parts of SGLD- and SGLD-FP-based trajectories before and after using control variates, and confidence intervals for estimation of  $f$ , are presented in Figure 6.

## Appendix A. Supplementary material for variance reduction for dependent sequences with applications to stochastic gradient MCMC.

**A.1. Proof of Proposition 2.1.** With notation  $\tilde{h} = h - \pi(h)$ , we can represent the variance of  $\pi_n(h)$ ,  $h \in \mathcal{H}$ , as

$$\mathbb{E} \left[ \left( \frac{1}{n} \sum_{k=0}^{n-1} \tilde{h}(X_k) \right)^2 \right] = \frac{1}{n^2} \left\{ \sum_{k=0}^{n-1} \mathbb{E}[\tilde{h}^2(X_k)] + 2 \sum_{\ell=1}^{n-1} \sum_{k=0}^{n-\ell-1} \mathbb{E}[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] \right\}.$$



**Figure 6.** Bayesian probabilistic matrix factorization from subsection 4.4. Left panel: test mean squared error (MSE) trajectory for SGLD (left) and SGLD-FP (right) with and without ESVM. Right panel: confidence intervals for test mean squared error trajectory for SGLD (left) and SGLD-FP (right).



Multiplying the both sides by  $n$  and subtracting  $\rho^{(h)}(0) + 2 \sum_{\ell=1}^{n-1} (1 - \ell n^{-1}) \rho^{(h)}(\ell)$ , we get

$$\begin{aligned} & n \mathbb{E} \left[ \left( \frac{1}{n} \sum_{k=0}^{n-1} \tilde{h}(X_k) \right)^2 \right] - \rho^{(h)}(0) - 2 \sum_{\ell=1}^{n-1} \left( 1 - \frac{\ell}{n} \right) \rho^{(h)}(\ell) \\ &= \frac{1}{n} \sum_{k=0}^{n-1} \left( \mathbb{E}[\tilde{h}^2(X_k)] - \rho^{(h)}(0) \right) + \frac{2}{n} \sum_{\ell=1}^{n-1} \sum_{k=0}^{n-\ell-1} \left( \mathbb{E}[\tilde{h}(X_k) \tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell) \right). \end{aligned}$$

It follows from the Cesàro mean theorem and (CS) that the right-hand side tends to zero as  $n \rightarrow \infty$ . Similarly,  $\rho^{(h)}(0) + 2 \sum_{\ell=0}^{n-1} (1 - \ell n^{-1}) \rho^{(h)}(\ell) \rightarrow \sum_{\ell \in \mathbb{Z}} \rho^{(h)}(\ell)$  as  $n \rightarrow \infty$ .

**A.2. Proof of Theorem 2.2.** Let us first start with a technical lemma the proof of which we postpone to the end of the section. In what follows, set  $\bar{V}_n(h) := \mathbb{E}[V_n(h)]$ .

**Lemma A.1.** *Let  $\mathcal{H}$  be a class of functions with constant mean and assume that (CS) and (CD) hold. Then, for any  $h \in \mathcal{H}$  and any  $h_1, h_2 \in \mathcal{H}$  with  $\|h_1 - h_2\|_{L^2(\pi)} \leq \varepsilon$ ,*

$$\begin{aligned} (1) \quad & |V_\infty(h) - \bar{V}_n(h)| \lesssim (R + \varsigma(1 - \lambda)^{-1}) b_n n^{-1} + \varsigma(1 - \lambda)^{-2} n^{-1} + \varsigma(1 - \lambda)^{-1} \lambda^{b_n/2}, \\ (2) \quad & |\bar{V}_n(h_1) - \bar{V}_n(h_2)| \lesssim \sqrt{RD} b_n n^{-1/2} + (R + \varsigma(1 - \lambda)^{-1}) b_n n^{-1} \\ & \quad + (\sqrt{R} n^{-1/2} + \sqrt{D}) b_n \varepsilon. \end{aligned}$$

Let  $h^*$  be a function in  $\mathcal{H}$  leading to the smallest  $\bar{V}_n(h)$ , that is,

$$h^* \in \arg \min_{h \in \mathcal{H}} \bar{V}_n(h).$$

For simplicity, we assume that  $h^*$  exists as all the following arguments can easily be adapted by considering an approximate minimizer. We decompose the excess of the asymptotic variance as

$$\begin{aligned} & V_\infty(\hat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty(h) \\ &= V_\infty(\hat{h}_{n,\varepsilon}) - \bar{V}_n(\hat{h}_{n,\varepsilon}) + \bar{V}_n(\hat{h}_{n,\varepsilon}) - \bar{V}_n(h^*) + \bar{V}_n(h^*) - \inf_{h \in \mathcal{H}} V_\infty(h) \\ (A.1) \quad & \leq 2 \sup_{h \in \mathcal{H}} |V_\infty(h) - \bar{V}_n(h)| + \bar{V}_n(\hat{h}_{n,\varepsilon}) - \bar{V}_n(h^*). \end{aligned}$$

To bound the first term in (A.1), we apply Lemma A.1 and obtain

$$\begin{aligned} (A.2) \quad & \sup_{h \in \mathcal{H}} |V_\infty(h) - \bar{V}_n(h)| \\ & \lesssim (R + \varsigma(1 - \lambda)^{-1}) b_n n^{-1} + \varsigma(1 - \lambda)^{-2} n^{-1} + \varsigma(1 - \lambda)^{-1} \lambda^{b_n/2}. \end{aligned}$$

It remains to bound the second term in (A.1). Let  $h_\varepsilon^* \in \mathcal{H}_\varepsilon$  be any closest to  $h^*$  point in  $L^2(\pi)$ -distance. By the definition of  $\hat{h}_{n,\varepsilon}$ ,  $V_n(\hat{h}_{n,\varepsilon}) - V_n(h_\varepsilon^*) \leq 0$ . Hence,

$$\begin{aligned} & \bar{V}_n(\hat{h}_{n,\varepsilon}) - \bar{V}_n(h^*) \\ & \leq \bar{V}_n(\hat{h}_{n,\varepsilon}) - \bar{V}_n(h^*) - (V_n(\hat{h}_{n,\varepsilon}) - V_n(h_\varepsilon^*)) \\ & = \bar{V}_n(\hat{h}_{n,\varepsilon}) - \bar{V}_n(h^*) - (V_n(\hat{h}_{n,\varepsilon}) - V_n(h^*)) + (V_n(h_\varepsilon^*) - V_n(h^*)) \\ (A.3) \quad & \leq \sup_{h \in \mathcal{H}_\varepsilon} \{ \bar{V}_n(h) - V_n(h) \} + (V_n(h^*) - \bar{V}_n(h^*)) + (V_n(h_\varepsilon^*) - V_n(h^*)). \end{aligned}$$

By assumption and the union bound, it holds for the first term in (A.3) that

$$\mathbb{P} \left( \sup_{h \in \mathcal{H}_\varepsilon} \{\bar{V}_n(h) - V_n(h)\} > t \right) \leq |\mathcal{H}_\varepsilon| \sup_{h \in \mathcal{H}_\varepsilon} \mathbb{P} \left( \bar{V}_n(h) - V_n(h) > t \right) \leq |\mathcal{H}_\varepsilon| \alpha_n(t).$$

The second term in (A.3) can be handled in the same way,

$$\mathbb{P}(V_n(h^*) - \bar{V}_n(h^*) > t) \leq \alpha_n(t).$$

The last term in (A.3) we represent as

$$V_n(h^*) - V_n(h_\varepsilon^*) = V_n(h^*) - V_n(h_\varepsilon^*) - (\bar{V}_n(h^*) - \bar{V}_n(h_\varepsilon^*)) + (\bar{V}_n(h^*) - \bar{V}_n(h_\varepsilon^*)).$$

Now the union bound implies

$$\mathbb{P} \left( V_n(h^*) - V_n(h_\varepsilon^*) - (\bar{V}_n(h^*) - \bar{V}_n(h_\varepsilon^*)) > t \right) \leq 2\alpha_n(t).$$

Furthermore, using Lemma A.1 and the fact that  $h_\varepsilon^*$  is  $\varepsilon$ -close to  $h^*$  in  $L^2(\pi)$ -distance,

$$\begin{aligned} & |\bar{V}_n(h^*) - \bar{V}_n(h_\varepsilon^*)| \\ & \lesssim \sqrt{RD}b_n n^{-1/2} + (R + \varsigma(1 - \lambda)^{-1})b_n n^{-1} + (\sqrt{R}n^{-1/2} + \sqrt{D})b_n \varepsilon. \end{aligned}$$

Combining these inequalities and substituting them into (A.3), we obtain, with probability at least  $1 - (|\mathcal{H}_\varepsilon| + 3)\alpha_n(t)$ ,

$$\begin{aligned} \text{(A.4)} \quad & \bar{V}_n(\hat{h}_{n,\varepsilon}) - \bar{V}_n(h^*) \\ & \lesssim t + \sqrt{RD}b_n n^{-1/2} + (R + \varsigma(1 - \lambda)^{-1})b_n n^{-1} + (\sqrt{R}n^{-1/2} + \sqrt{D})b_n \varepsilon. \end{aligned}$$

Substituting (A.2) and (A.4) into (A.1) we conclude that, with the same probability,

$$\begin{aligned} V_\infty(\hat{h}_{n,\varepsilon}) - \inf_{h \in \mathcal{H}} V_\infty(h) & \lesssim t + (\sqrt{R}n^{-1/2} + \sqrt{D})b_n \varepsilon + \sqrt{RD}b_n n^{-1/2} \\ & \quad + (R + \varsigma(1 - \lambda)^{-1})b_n n^{-1} + \varsigma(1 - \lambda)^{-2}n^{-1} + \varsigma(1 - \lambda)^{-1}\lambda^{b_n/2}, \end{aligned}$$

where we have used the notation of Theorem 2.2. The proof is completed by taking  $t = \alpha_n^{-1}(\delta/2|\mathcal{H}_\varepsilon|)$  and assuming that  $|\mathcal{H}_\varepsilon| \geq 3$  (this involves no loss of generality). We are left with the task of proving Lemma A.1.

*Proof of Lemma A.1.* Let us first find a leading term in sample autocovariance function. Recall that for any  $h \in \mathcal{H}$ ,  $\tilde{h} = h - \pi(h)$ . By expanding the brackets and adding/subtracting  $\pi(h)$  in the definition (2.3), we get, for any  $|\ell| \leq b_n$ ,

$$\text{(A.5)} \quad \rho_n^{(h)}(\ell) = A_{n,1}^{(h)}(\ell) + A_{n,2}^{(h)}(\ell) + A_{n,3}^{(h)}(\ell),$$

where, for  $0 \leq \ell \leq b_n$ ,

$$\begin{aligned} A_{n,1}^{(h)}(\ell) &= A_{n,1}^{(h)}(-\ell) := n^{-1} \sum_{k=0}^{n-\ell-1} \tilde{h}(X_k) \tilde{h}(X_{k+\ell}), \\ A_{n,2}^{(h)}(\ell) &= A_{n,2}^{(h)}(-\ell) := -n^{-1} \pi_n(\tilde{h}) \left\{ \sum_{k=0}^{n-\ell-1} \tilde{h}(X_k) + \sum_{k=\ell}^{n-1} \tilde{h}(X_k) \right\}, \\ A_{n,3}^{(h)}(\ell) &= A_{n,3}^{(h)}(-\ell) := (1 - \ell/n) \pi_n^2(\tilde{h}). \end{aligned}$$

It follows that the leading term in this decomposition is  $A_{n,1}^{(h)}(\ell)$ . The remainder terms  $A_{n,2}^{(h)}(\ell)$  and  $A_{n,3}^{(h)}(\ell)$  can be bounded, under assumptions (CS) and (CD), as follows:

$$\begin{aligned} \left| \mathbb{E}[A_{n,3}^{(h)}(\ell)] \right| &\leq n^{-2} \left\{ \sum_{k=0}^{n-1} \mathbb{E}[\tilde{h}^2(X_k)] + 2 \sum_{\ell=0}^{n-1} \sum_{k=1}^{n-\ell-1} \mathbb{E}[\tilde{h}(X_k) \tilde{h}(X_{k+\ell})] \right\} \\ &\leq n^{-1} \left\{ \rho^{(h)}(0) + 2 \sum_{\ell=0}^{n-1} (1 - \ell n^{-1}) \rho^{(h)}(\ell) \right\} + 3Rn^{-1} \\ \text{(A.6)} \qquad \qquad \qquad &\leq Cn^{-1}, \end{aligned}$$

where  $C := 2\zeta(1 - \lambda)^{-1} + 3R$ . In the same manner we conclude that

$$\begin{aligned} \left| \mathbb{E}[A_{n,2}^{(h)}(\ell)] \right| &\leq n^{-1} \mathbb{E}^{1/2}[\pi_n^2(\tilde{h})] \left\{ \mathbb{E}^{1/2} \left[ \left( \sum_{k=0}^{n-\ell-1} \tilde{h}(X_k) \right)^2 + \left( \sum_{k=\ell}^{n-1} \tilde{h}(X_k) \right)^2 \right] \right\} \\ \text{(A.7)} \qquad \qquad \qquad &\leq 2Cn^{-1}. \end{aligned}$$

The last two bounds show that the last two terms in (A.5) are of order  $n^{-1}$ . Having disposed of this preliminary step, we can now return to statements of the lemma. ■

*Statement 1.* From decomposition (A.5) and bounds (A.6), (A.7), we deduce that

$$\begin{aligned} |\bar{V}_n(h_1) - \bar{V}_n(h_2)| &= \sum_{|\ell| \leq b_n} w_n(\ell) \mathbb{E} \left[ \rho_n^{(h_1)}(\ell) - \rho_n^{(h_2)}(\ell) \right] \\ &\leq 2b_n \max_{|\ell| \leq b_n} \left| \mathbb{E} \left[ A_{n,1}^{(h_1)}(\ell) - A_{n,1}^{(h_2)}(\ell) \right] \right| + 12Cb_n n^{-1}. \end{aligned}$$

With notation  $\tilde{h}_{12} = \tilde{h}_1 - \tilde{h}_2$ , it follows, for any  $0 \leq \ell \leq b_n$ , that

$$A_{n,1}^{(h_1)}(\ell) - A_{n,1}^{(h_2)}(\ell) = n^{-1} \sum_{k=0}^{n-\ell-1} \left( \tilde{h}_1(X_k) \tilde{h}_{12}(X_{k+\ell}) + \tilde{h}_{12}(X_k) \tilde{h}_2(X_{k+\ell}) \right).$$

Using the Cauchy–Schwarz inequality (twice) and (CS), we have

$$\begin{aligned} \mathbb{E} \left[ \sum_{k=0}^{n-\ell-1} \tilde{h}_1(X_k) \tilde{h}_{12}(X_{k+\ell}) \right] &\leq \mathbb{E}^{1/2} \left[ \sum_{k=0}^{n-\ell-1} \tilde{h}_1^2(X_k) \right] \mathbb{E}^{1/2} \left[ \sum_{k=0}^{n-\ell-1} \tilde{h}_{12}^2(X_{k+\ell}) \right] \\ &\leq \sqrt{R + (n - \ell) \rho^{(\tilde{h}_1)}(0)} \sqrt{R + (n - \ell) \rho^{(\tilde{h}_{12})}(0)}. \end{aligned}$$

We now apply this argument again and obtain

$$\mathbb{E}\left[A_{n,1}^{(h_1)}(\ell) - A_{n,1}^{(h_2)}(\ell)\right] \lesssim Rn^{-1} + \sqrt{RD}n^{-1/2} + (\sqrt{R}n^{-1/2} + \sqrt{D})\|\tilde{h}_1 - \tilde{h}_2\|_{L^2(\pi)}.$$

Finally, since  $\|\tilde{h}_1 - \tilde{h}_2\|_{L^2(\pi)} \leq 2\|h_1 - h_2\|_{L^2(\pi)}$ , we conclude

$$\begin{aligned} & |\bar{V}_n(h_1) - \bar{V}_n(h_2)| \\ & \lesssim (R + C)b_n n^{-1} + \sqrt{RD}b_n n^{-1/2} + (\sqrt{R}n^{-1/2} + \sqrt{D})b_n \|h_1 - h_2\|_{L^2(\pi)}. \end{aligned}$$

*Statement 2.* Let us denote

$$V_{n,\rho}(h) = \sum_{|\ell| \leq b_n} w_n(\ell) \rho^{(h)}(\ell).$$

With this notation, we have the following decomposition:

$$(A.8) \quad |V_\infty(h) - \bar{V}_n(h)| \leq |V_\infty(h) - V_{n,\rho}(h)| + |V_{n,\rho}(h) - \bar{V}_n(h)|.$$

To bound the first term in the right-hand side of (A.8), we represent it as

$$|V_{n,\rho}(h) - V_\infty(h)| \leq \sum_{|\ell| \leq b_n} |1 - w_n(\ell)| |\rho^{(h)}(\ell)| + \sum_{|\ell| > b_n} |\rho^{(h)}(\ell)|.$$

Using (CD) and the fact that  $w_n(\ell) = 1$  for  $\ell \in [-b_n/2, b_n/2]$ , we obtain

$$\sum_{|\ell| \leq b_n} |1 - w_n(\ell)| |\rho^{(h)}(\ell)| = 2 \sum_{\ell=b_n/2}^{b_n} |1 - w_n(\ell)| |\rho^{(h)}(\ell)| \leq 2\zeta(1 - \lambda)^{-1} \lambda^{b_n/2}.$$

In the same manner we can see that

$$\sum_{|s| > b_n} |\rho^{(h)}(s)| \leq 2\zeta(1 - \lambda)^{-1} \lambda^{b_n}.$$

Combining the last two bounds we conclude that

$$(A.9) \quad |V_{n,\rho}(h) - V_\infty(h)| \leq 4\zeta(1 - \lambda)^{-1} \lambda^{b_n/2}.$$

Now let us turn to the second term in the right-hand side of (A.8). The decomposition (A.5) and the bounds (A.6), (A.7) yield

$$\begin{aligned} |\bar{V}_n(h) - V_{n,\rho}(h)| &= \sum_{|\ell| \leq b_n} w_n(\ell) \left( \mathbb{E}[\rho_n^{(h)}(\ell)] - \rho^{(h)}(\ell) \right) \\ &\leq \sum_{|\ell| \leq b_n} \left| \mathbb{E}[A_{n,1}^{(h)}(\ell)] - \rho^{(h)}(\ell) \right| + 6Cb_n n^{-1}. \end{aligned}$$

Using (CS) and (CD), it follows that

$$\begin{aligned} \left| \mathbb{E}[A_{n,1}^{(h)}(\ell)] - \rho^{(h)}(\ell) \right| &\leq n^{-1} \sum_{k=0}^{n-\ell-1} \left| \mathbb{E}[\tilde{h}(X_k) \tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell) \right| + \ell n^{-1} \rho^{(h)}(\ell) \\ &\leq Rn^{-1} + \varsigma \ell \lambda^\ell n^{-1}. \end{aligned}$$

Combining these, we get

$$(A.10) \quad |\bar{V}_n(h) - V_{n,\rho}(h)| \leq (12\zeta(1 - \lambda)^{-1} + 20R)b_n n^{-1} + 2\zeta(1 - \lambda)^{-2} n^{-1}.$$

Finally, we obtain the desired conclusion by substituting (A.9) and (A.10) into (A.8).

**A.3. Proof of Proposition 3.2.** 1. The sequence  $(\rho^{(h)}(\ell))_{\ell=0}^\infty$  is symmetric and positive semidefinite by construction. By the Markov property, for any  $k, \ell \in \mathbb{N}$ ,

$$E_\xi[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell) = \bar{E}_\zeta[\tilde{h}(X)\phi_\ell(X) - \tilde{h}(X')\phi_\ell(X')],$$

where we denote  $\phi_\ell(x) := P^\ell \tilde{h}(x)$  and  $\zeta \in \Pi(\xi P^k, \pi)$  is the optimal coupling of  $\xi P^k$  and  $\pi$  in  $W_2^d$ -distance,  $(X, X') \sim \zeta$ . Note that

$$\begin{aligned} |\bar{E}_\zeta[\tilde{h}(X)\phi_\ell(X) - \tilde{h}(X')\phi_\ell(X')]| &\leq \{\bar{E}_\zeta[\{\tilde{h}(X) - \tilde{h}(X')\}^2]\}^{1/2} \{\bar{E}_\zeta[\{\phi_\ell(X')\}^2]\}^{1/2} \\ &\quad + \{\bar{E}_\zeta[\{\phi_\ell(X) - \phi_\ell(X')\}^2]\}^{1/2} \{\bar{E}_\zeta[\{\tilde{h}(X)\}^2]\}^{1/2}. \end{aligned}$$

It is easy to check that  $\phi_\ell$  is a Lipschitz function,

$$|\phi_\ell(x) - \phi_\ell(x')| \leq L W_2^d(\delta_x P^\ell, \delta_{x'} P^\ell) \leq L \Delta_2^\ell d(x, x').$$

Since the Markov kernel  $P$  is  $W_2$ -geometrically ergodic, we get

$$\begin{aligned} |\bar{E}_\zeta[\tilde{h}(X)\phi_\ell(X) - \tilde{h}(X')\phi_\ell(X')]| \\ \leq L W_2^d(\xi P^k, \pi) \{\bar{E}_\zeta[\{\phi_\ell(X')\}^2]\}^{1/2} + L \Delta_2^\ell W_2^d(\xi P^k, \pi) \{\bar{E}_\zeta[\{\tilde{h}(X)\}^2]\}^{1/2}. \end{aligned}$$

Let us compute  $\bar{E}_\zeta[\{\phi_\ell(X')\}^2] = \pi(\phi_\ell^2)$ . Since  $P$  is  $W_2^d$ -geometrically ergodic, we have  $W_2^d(\delta_y P^\ell, \pi) \leq \Delta_2^\ell W_2^d(\delta_y, \pi)$ . Note also that  $\pi(\tilde{h}) = 0$  implies  $\pi(\phi_\ell) = 0$ ; hence

$$(A.11) \quad \pi(\phi_\ell^2) = \int \left[ \phi_\ell(y) - \int \phi_\ell(x) \pi(dx) \right]^2 \pi(dy) \leq L^2 \Delta_2^{2\ell} \int \{W_2^d(\delta_y, \pi)\}^2 \pi(dy).$$

Finally, we need to compute  $\bar{E}_\zeta[\{\tilde{h}(X)\}^2] = \xi P^k(\tilde{h}^2)$ . For an arbitrary  $\hat{x} \in \mathcal{X}$ ,

$$|\tilde{h}(x)|^2 = \left| \int \{h(x) - h(y)\} \pi(dy) \right|^2 \leq 2L^2 \left( d^2(x, \hat{x}) + \int d^2(x, \hat{x}) \pi(dx) \right).$$

In order to bound  $\xi P^k(d^2(x, \hat{x}))$ , we write

$$\begin{aligned} \int d^2(x, \hat{x}) \xi P^k(dx) &= \iint d^2(x, \hat{x}) \zeta(dx dx') \leq 2 \iint d^2(x, x') \zeta(dx dx') \\ &\quad + 2 \int d^2(x, \hat{x}) \pi(dx) \leq 2 \Delta_2^{2k} \{W_2^d(\xi, \pi)\}^2 + 2 \int d^2(x, \hat{x}) \pi(dx). \end{aligned}$$

Hence

$$\xi P^k(\tilde{h}^2) \leq 4L^2 \int d^2(x, \hat{x}) \pi(dx) + 2L^2 \Delta_2^{2k} \{W_2^d(\xi, \pi)\}^2,$$

and  $|E_\xi[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell)| \leq A_1 L^2 \Delta_2^{k+\ell} W_2^d(\xi, \pi)$ , where

$$(A.12) \quad A_1 := 2 \inf_{\hat{x} \in \mathcal{X}} W_2^d(\delta_{\hat{x}}, \pi) + 2W_2^d(\xi, \pi) + \left[ \int \{W_2^d(\delta_y, \pi)\}^2 \pi(dy) \right]^{1/2}.$$

Summing the last inequality with respect to  $k$ , we obtain

$$(A.13) \quad \sum_{k=0}^{\infty} |\mathbb{E}_{\xi}[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell)| \leq A_1 L^2 (1 - \Delta_2)^{-1} \Delta_2^{\ell} W_2(\xi, \pi).$$

Hence, the second assumption in (CS) holds with  $R$  defined in (3.2). The third assertion clearly follows from (A.13). To check (CD), we write

$$(A.14) \quad \begin{aligned} |\rho^{(h)}(\ell)| &= \left| \int \tilde{h}(x) [\delta_x P^{\ell}(h) - \pi(h)] \pi(dx) \right| \leq L \int |\tilde{h}(x)| W_2^d(\delta_x P^{\ell}, \pi) \pi(dx) \\ &\leq L \Delta_2^{\ell} \int |\tilde{h}(x)| W_2^d(\delta_x, \pi) \pi(dx) \leq L \Delta_2^{\ell} \sqrt{D} \left[ \int \{W_2^d(\delta_x, \pi)\}^2 \pi(dx) \right]^{1/2}. \end{aligned}$$

Hence (CD) holds with  $\lambda = \Delta_2$  and  $\varsigma = L\sqrt{D}[\int \{W_2^d(\delta_x, \pi)\}^2 \pi(dx)]^{1/2}$ .

2. The proof essentially relies on [17]. Denote  $Z_n(h) := (h(X_0), \dots, h(X_{n-1}))$ , and recall the representation (2.2). It follows from [5, section 5.2] that  $V_n(h)$  can be represented as a quadratic form

$$V_n(h) = \langle A_n Z_n(h), Z_n(h) \rangle,$$

where  $A_n = n^{-1}(\mathbf{I} - n^{-1}E)W(\mathbf{I} - n^{-1}E)$ ,  $E$  is an  $n \times n$  matrix with elements  $E_{j,k} = 1$  for any  $1 \leq j, k \leq n$ , and  $W$  is a Toeplitz matrix with elements  $W_{j,k} = w_n(j - k)$ . Note that  $V_n(h)$  is invariant to shifts and, in particular,  $V_n(h) = V_n(\tilde{h})$ . It is straightforward to show that  $\|A_n\| \leq 2bn^{-1}$ ; see [5, Lemma 9]. Furthermore, [5, Corollary 18] implies

$$(A.15) \quad \mathbb{P}_{\xi}(|V_n(h) - \mathbb{E}_{\xi}[V_n(h)]| \geq t) \leq 2 \exp\left(-\frac{(1 - \Delta_2)^2 t^2}{c\alpha L^2 (\mathbb{E}_{\xi}[\|A_n Z_n(\tilde{h})\|^2] + t\|A_n\|)}\right),$$

where  $c > 0$  is some universal constant. By the Cauchy–Schwarz inequality,  $\|A_n Z_n(h)\|^2 \leq \|A_n\|^2 \|Z_n(h)\|^2$ . Moreover, using (CS), we get

$$(A.16) \quad \mathbb{E}_{\xi}[\|Z_n(\tilde{h})\|^2] \leq R + n \text{Var}_{\pi}[h] \leq R + n \sup_{h \in \mathcal{H}} \text{Var}_{\pi}[h] = R + nD.$$

The statement follows from substitution (A.16) into (A.15).

**A.4. Proof of Proposition 3.3.** 1. Proceeding similarly to Appendix A.3, we use the Markov property to write, for  $k, \ell \in \mathbb{N}$ ,

$$\begin{aligned} |\mathbb{E}_{\xi}[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell)| &\leq \{\bar{\mathbb{E}}_{\zeta}[\{\tilde{h}(X) - \tilde{h}(X')\}^2]\}^{1/2} \{\bar{\mathbb{E}}_{\zeta}[\{\phi_{\ell}(X')\}^2]\}^{1/2} \\ &\quad + \{\bar{\mathbb{E}}_{\zeta}[\{\phi_{\ell}(X) - \phi_{\ell}(X')\}^2]\}^{1/2} \{\bar{\mathbb{E}}_{\zeta}[\{\tilde{h}(X)\}^2]\}^{1/2}, \end{aligned}$$

where  $\zeta \in \Pi(\xi P^k, \pi)$  is the optimal coupling of  $\xi P^k$  and  $\pi$  in  $W_1^d$  distance,  $(X, X') \sim \zeta$ . Since function  $\tilde{h}$  is bounded and Lipschitz,

$$\bar{\mathbb{E}}_{\zeta}[\{\tilde{h}(X) - \tilde{h}(X')\}^2] \leq 2B\bar{\mathbb{E}}_{\zeta}[\|\tilde{h}(X) - \tilde{h}(X')\|] \leq 2LB\Delta_1^k W_1^d(\xi, \pi).$$

Similarly, using that  $\phi_{\ell}$  is bounded and Lipschitz (see Appendix A.3 for the details),

$$\bar{\mathbb{E}}_{\zeta}[\{\phi_{\ell}(X) - \phi_{\ell}(X')\}^2] \leq 2LB\Delta_1^{k+\ell} W_1^d(\xi, \pi).$$

Proceeding as in (A.11), we obtain

$$\bar{E}_\zeta[\{\phi_\ell(X')\}^2] \leq \iint [\phi_\ell(x) - \phi_\ell(y)]^2 \pi(dx)\pi(dy) \leq 2BL\Delta_1^\ell \int W_1^d(\delta_y, \pi) \pi(dy).$$

Using the simple bound  $\bar{E}_\zeta[\{\tilde{h}(X)\}^2] \leq 4B^2$ , it holds that

$$(A.17) \quad |\mathbb{E}_\xi[\tilde{h}(X_k)\tilde{h}(X_{k+\ell})] - \rho^{(h)}(\ell)| \leq 2BC'_2\Delta_1^{(k+\ell)/2}$$

with

$$(A.18) \quad C'_2 = 2L\{W_1^d(\xi, \pi)\}^{1/2} \left\{ \int W_1^d(\delta_x, \pi)\pi(dx) \right\}^{1/2} + \{2LBW_1^d(\xi, \pi)\}^{1/2}.$$

Hence, the second assumption in (CS) holds with

$$(A.19) \quad R = 2BC'_2(1 - \Delta_1^{1/2})^{-1},$$

and the third one follows from (A.17). Proceeding as in (A.14),

$$|\rho^{(h)}(\ell)| \leq 2LB\Delta_1^\ell \int W_1^d(\delta_x, \pi)\pi(dx),$$

and (CD) holds with  $\lambda = \Delta_1$  and

$$(A.20) \quad \varsigma = 2LB \int W_1^d(\delta_x, \pi)\pi(dx).$$

2. Without loss of generality, we assume that  $\|h\|_\infty \leq 1$ . By Minkowski's inequality,

$$\|V_n(h) - \mathbb{E}_\xi V_n(h)\|_{\xi,p} \leq \sum_{\ell=-b_n+1}^{b_n} w_n(\ell) \|\rho_n^{(h)}(\ell) - \mathbb{E}_\xi[\rho_n^{(h)}(\ell)]\|_{\xi,p},$$

where  $\|\cdot\|_{\xi,p} := (\mathbb{E}_\xi[\cdot]^p)^{1/p}$ . For  $\ell \in \mathbb{N}_0$ , we get

$$\begin{aligned} \rho_n^{(h)}(\ell) &= \frac{1}{n} \sum_{k=0}^{n-\ell-1} (h(X_k) - \pi_n(h))(h(X_{k+\ell}) - \pi_n(h)) = \frac{1}{n} \sum_{k=0}^{n-\ell-1} \tilde{h}(X_k)\tilde{h}(X_{k+\ell}) \\ &\quad - \frac{1}{n} (\pi_n(h) - \pi(h))^2 + \frac{\pi_n(h) - \pi(h)}{n} \sum_{k=0}^{\ell-1} [\tilde{h}(X_k) + \tilde{h}(X_{n-\ell+k})] \\ &=: T_1 + T_2 + T_3. \end{aligned}$$

Hence,

$$\rho_n^{(h)}(\ell) - \mathbb{E}_\xi[\rho_n^{(h)}(\ell)] = (T_1 - \mathbb{E}_\xi[T_1]) + (T_2 - \mathbb{E}_\xi[T_2]) + (T_3 - \mathbb{E}_\xi[T_3]) =: \bar{T}_1 + \bar{T}_2 + \bar{T}_3.$$

Now we proceed with estimating  $\|\bar{T}_1\|_{\xi,p}$ . By [19, Theorem 2] and Lemma A.2, setting  $\beta_{r,\ell} = 1 \wedge \Delta_1^{r-\ell}$  and  $A_1 = 4 \vee 256L\{W_1(\xi, \pi) + 2 \inf_{\hat{x} \in X} W_1(\delta_{\hat{x}}, \pi)\}$ ,

$$(A.21) \quad \|\bar{T}_1\|_{\xi,p}^p \leq \frac{(2p-2)!}{(p-1)!} e^p \left[ \left( nA_1 \sum_{r=0}^{n-1} \beta_{r,\ell} \right)^{p/2} \vee nA_1 16^{p-2} \sum_{r=0}^{n-1} (r+1)^{p-2} \beta_{r,\ell} \right].$$

It can be easily seen that

$$(A.22) \quad \sum_{r=0}^{n-1} (r+1)^{p-1} \beta_{r,\ell} \leq \frac{\ell^{p-1}}{p-1} + \frac{p!}{\Delta_1^2} \left[ \frac{1}{\log^p(1/\Delta_1)} + \ell^p \right], \quad \sum_{r=0}^{n-1} \beta_{r,\ell} \leq \frac{2\ell}{(1-\Delta_1)}.$$

Substituting (A.22) into (A.21) and using Stirling's formula,

$$\|\bar{T}_1\|_{\xi,p}^p \leq 2^{2p} p^p \left[ \left( \frac{2A_1 \ell n}{1-\Delta_1} \right)^{p/2} + nA_1 16^p \left\{ \ell^p + \frac{2^p p^p p^{1/2}}{e^p \Delta_1^2} \left( \frac{1}{\log^p(1/\Delta_1)} + \ell^p \right) \right\} \right].$$

Since  $\ell \leq b_n$ , we obtain the following final bound on  $\bar{T}_1$ ,

$$(A.23) \quad \|\bar{T}_1\|_{\xi,p} \leq 4p \left[ \left( \frac{2b_n A_1}{n(1-\Delta_1)} \right)^{1/2} + 32b_n n^{1/p-1} A_1^{1/p} \left( 1 + \frac{p(1+\log(1/\Delta_1))}{e\Delta_1^{2/p} \log(1/\Delta_1)} \right) \right].$$

Let us consider now  $\bar{T}_2$  and  $\bar{T}_3$ . Using (WE),

$$\|\pi_n(h) - \pi(h)\|_{\xi,p} \leq n^{-1} \left\| \sum_{k=0}^{n-1} h(X_k) - \mathbb{E}_\xi[h(X_k)] \right\|_{\xi,p} + \frac{LW_1(\xi, \pi)}{n(1-\Delta)}.$$

By Lemma A.3 and [19, Theorem 2], setting  $A_2 = 4L\{W_1(\xi, \pi) + 2 \inf_{\hat{x} \in X} W_1(\delta_{\hat{x}}, \pi)\}$ ,

$$\left\| \sum_{k=0}^{n-1} h(X_k) - \mathbb{E}_\xi[h(X_k)] \right\|_{\xi,p} \leq 4p \left[ \frac{n^{1/2} A_2^{1/2}}{\sqrt{1-\Delta_1}} \vee \frac{2pn^{1/p} A_2^{1/p}}{\Delta_1^{1/p} e \log(1/\Delta_1)} \right].$$

Now it holds for  $\bar{T}_2$ ,

$$(A.24) \quad \begin{aligned} \|\bar{T}_2\|_{\xi,p} &\leq 2 \|\pi_n(h) - \pi(h)\|_{\xi,2p}^2 \leq \frac{4}{n^2} \left\| \sum_{k=0}^{n-1} h(X_k) - \mathbb{E}_\xi[h(X_k)] \right\|_{\xi,2p}^2 + \frac{6L^2 W_1^2(\xi, \pi)}{n^2(1-\Delta)^2} \\ &\leq 2^6 p^2 \left[ \frac{A_2}{n(1-\Delta_1)} \vee \frac{4p^2 n^{2/p} A_2^{2/p}}{n^2 \Delta_1^{2/p} e^2 \log^2(1/\Delta_1)} \right] + \frac{4L^2 W_1^2(\xi, \pi)}{n^2(1-\Delta_1)^2}. \end{aligned}$$

Finally, since  $\ell \leq b_n$  and  $h$  is bounded,

$$(A.25) \quad \|\bar{T}_3\|_{\xi,p} \leq 16b_n n^{-1}.$$

Using (A.23), (A.24), and (A.25), we get

$$(A.26) \quad \|V_n(h) - \mathbb{E}_\xi[V_n(h)]\|_{\xi,p} \leq 2b_n \left[ \frac{C'_1 p b_n^{1/2}}{n^{1/2}} + \frac{C'_2 p^2 b_n}{n^{1-1/p}} + \frac{C'_3 p^4}{n^{2-2/p}} \right],$$



where

$$(A.27) \quad \begin{aligned} C'_1 &= \frac{4(2A_1)^{1/2}}{\sqrt{1-\Delta_1}}, \quad C'_2 = \frac{2^7 A_1^{1/p}(1+2\log(1/\Delta_1))}{\Delta_1^{2/p} \log(1/\Delta_1)} + \frac{2^6 A_2}{1-\Delta_1} + 12, \\ C'_3 &= \frac{2^8 A_2^{2/p}}{e^2 \Delta_1^{2/p} \log^2(1/\Delta_1)} + \frac{A_2^2}{4(1-\Delta_1)^2}. \end{aligned}$$

Under the assumption  $p < n^{1/2}$ , we obtain

$$(A.28) \quad \|V_n(h) - E_\xi[V_n(h)]\|_{\xi,p} \leq b_n \|h\|_\infty^2 \left[ \frac{C_{R,1} p b_n^{1/2}}{n^{1/2}} + \frac{C_{R,2} p^2 b_n}{n^{1-1/p}} \right]$$

with  $C_{R,1} = 2C'_1$ ,  $C_{R,2} = 2C'_2 + 2C'_3$ . Now the statement follows from Markov's inequality.

**Lemma A.2.** Assume (WE)-1. Let  $h \in \text{Lip}_{b,d}(L, B)$ . For  $i, m \in \mathbb{N}_0$ , we define  $\tilde{g}_{i,m}(x, x') = \tilde{h}(x)\tilde{h}(x') - c_{\xi,i,m}$ , where  $c_{\xi,i,m} := E_\xi[g_{i,m}(X_i, X_{i+m})]$ . For  $p, r, m \in \mathbb{N}_0$ , let

$$C_{p,r,m}^{(h)} := \sup \left| \text{cov}_\xi \left( \prod_{k=1}^u \tilde{g}_{i_k,m}(X_{i_k}, X_{i_k+m}), \prod_{k=1}^v \tilde{g}_{j_k,m}(X_{j_k}, X_{j_k+m}) \right) \right|,$$

where the supremum is taken over all  $0 \leq i_1 \leq \dots \leq i_u < i_u + r \leq j_1 \leq \dots \leq j_v \leq n$  with  $u + v = p$ . Then, for any  $p, r \in \mathbb{N}$ ,

$$C_{p,r,m}^{(h)} \leq \begin{cases} 2^{2p+2} B^{2p}, & r \leq m, \\ L \{W_1(\xi, \pi) + 2 \inf_{\hat{x} \in X} W_1(\delta_{\hat{x}}, \pi)\} v 2^{4p} B^{2p-1} \Delta_1^{r-m}, & r > m. \end{cases}$$

*Proof.* Define the function

$$G_{i_1, \dots, i_u, m}(x_{i_1}, x_{i_1+m}, \dots, x_{i_u}, x_{i_u+m}) := \prod_{k=1}^u \tilde{g}_{i_k, m}(x_{i_k}, x_{i_k+m}).$$

Let  $D_{i_1, \dots, i_u, m} = G_{i_1, \dots, i_u, m}(X_{i_1}, X_{i_1+m}, \dots, X_{i_u}, X_{i_u+m})$ . Since  $\|G_{i_1, \dots, i_u, m}\|_\infty \leq (2B)^{2u}$  and  $\|G_{j_1, \dots, j_v, m}\|_\infty \leq (2B)^{2v}$ , we get  $C_{p,r,m}^{(h)} \leq 2^{2p+2} B^{2p}$ . Now let  $m < r$ . Using Markov's property,

$$\begin{aligned} \text{cov}_\xi(D_{i_1, \dots, i_u, m}, D_{j_1, \dots, j_v, m}) &= E_\xi[(D_{i_1, \dots, i_u, m} - E_\xi[D_{i_1, \dots, i_u, m}]) (P^{j_1-i_u-m} \varphi(X_{i_u+m}) - \pi(\varphi))], \end{aligned}$$

where

$$\varphi(x) := E_x[G_{j_1, \dots, j_v, m}(x, X_m, X_{j_2-j_1}, X_{j_2-j_1+m}, \dots, X_{j_v-j_1}, X_{j_v-j_1+m})].$$

It follows from Lemma A.4 that  $\|\varphi\|_{\text{Lip}} \leq Lv 2^{4v-1} B^{2v-1}$ . By [18, Theorem 20.1.2]  $\|P^{j_1-i_u-m} \varphi\|_{\text{Lip}} \leq \Delta_1^{j_1-i_u-m} \|\varphi\|_{\text{Lip}}$ , and hence

$$|P^{j_1-i_u-m} \varphi(x) - \pi(\varphi)| \leq Lv 2^{4v-1} B^{2v-1} \Delta_1^{j_1-i_u-m} W_1(\delta_x, \pi).$$

This yields

$$|\text{cov}_\xi(D_{i_1, \dots, i_u, m}, D_{j_1, \dots, j_v, m})| \leq Lp 2^{4p} B^{2p-1} \Delta_1^{j_1-i_u-m} E_\xi[W_1(\delta_{X_{i_u+m}}, \pi)].$$

For a fixed  $\hat{x} \in \mathbf{X}$ , by the triangle inequality,

$$W_1(\delta_x, \pi) \leq W_1(\delta_x, \delta_{\hat{x}}) + W_1(\delta_{\hat{x}}, \pi) = d(x, \hat{x}) + W_1(\delta_{\hat{x}}, \pi).$$

Since  $\mathbb{E}_\xi[d(X_{i_u+m}, \hat{x})] \leq W_1(\delta_\xi P^{i_u+m}, \delta_{\hat{x}})$ , we get

$$\mathbb{E}_\xi[d(X_{i_u+m}, \hat{x})] \leq W_1(\xi P^{i_u+m}, \pi) + W_1(\delta_{\hat{x}}, \pi) \leq \Delta^{i_u+m} W_1(\xi, \pi) + W_1(\delta_{\hat{x}}, \pi),$$

showing that  $\mathbb{E}_\xi[W_1(\delta_{X_{i_u+m}}, \pi)] \leq W_1(\xi, \pi) + 2W_1(\delta_{\hat{x}}, \pi)$ . The proof is complete.  $\blacksquare$

**Lemma A.3.** Assume (WE)-1. Let  $h \in \text{Lip}_{b,d}(L, B)$ . For  $p, r \in \mathbb{N}_0$ , we define

$$C_{p,r} := \sup \left| \text{cov}_\xi \left( \prod_{k=1}^u (h(X_{i_k}) - \mathbb{E}_\xi[h(X_{i_k})]), \prod_{k=1}^v (h(X_{j_k}) - \mathbb{E}_\xi[h(X_{j_k})]) \right) \right|,$$

where the supremum is taken over all  $0 \leq i_1 \leq \dots \leq i_u < i_u + r \leq j_1 \leq \dots \leq j_v \leq n$  with  $u + v = p$ . Then for any  $p, r \in \mathbb{N}$ ,

$$C_{p,r} \leq L \{W_1(\xi, \pi) + 2 \inf_{\hat{x} \in \mathbf{X}} W_1(\delta_{\hat{x}}, \pi)\} p 2^{2p} B^{2p-1} \Delta_1^r.$$

*Proof.* The proof is along the same lines as Lemma A.2 and is omitted.  $\blacksquare$

**Lemma A.4.** Assume (WE)-1. Set

$$\varphi(x) = \mathbb{E}_x[G_{j_1, \dots, j_v, m}(x, X_m, X_{j_2-j_1}, X_{j_2-j_1+m}, \dots, X_{j_v-j_1}, X_{j_v-j_1+m})],$$

where  $G_{j_1, \dots, j_v, m}$  defined in Lemma A.2. Then

$$\|\varphi\|_{\text{Lip}} \leq Lv 2^{4v-1} B^{2v-1}.$$

*Proof.* We split the proof into two parts. First, we estimate Lipschitz constant of  $g(x) = \mathbb{E}_x[\prod_{k=1}^v \tilde{h}(X_{i_k}) \tilde{h}(X_{i_k+m})]$  for  $0 = i_1 \leq i_2 \leq \dots \leq i_v$  and  $m > 0$ . Note that  $g(x) = \tilde{h}^{n_1}(x) \mathbb{E}_x[\prod_{k=2}^b \tilde{h}^{n_k}(X_{m_k})]$ , where  $0 = m_1 < m_2 < \dots < m_b$  are distinct indices among  $(i_1, i_1 + m, \dots, i_v, i_v + m)$  and  $(n_1, \dots, n_b)$  are their associated multiplicities ( $\sum_{k=1}^b n_k = 2v$ ). Hence, applying Lemma A.5 with  $f_i = \tilde{h}$  and  $K = 2B$ , we get  $\|g\|_{\text{Lip}} \leq 2Lv(2B)^{2v-1}$ . Now we estimate the Lipschitz constant of

$$(A.29) \quad \varphi(x) = \mathbb{E}_x \left[ \prod_{k=1}^v (\tilde{h}(X_{j_k-j_1}) \tilde{h}(X_{j_k-j_1+m}) - c_{\xi, j_k, m}) \right],$$

where  $c_{\xi, j_k, m} := \mathbb{E}_\xi[\tilde{h}(X_{j_k}) \tilde{h}(X_{j_k+m})]$ . Expanding (A.29), we obtain

$$(A.30) \quad \varphi(x) = \sum_{(\delta_1, \dots, \delta_v)} (-1)^{\sum_k \delta_k} \mathbb{E}_x \left[ \prod_{k=1}^v \tilde{h}^{\delta_k}(X_{j_k}) \tilde{h}^{\delta_k}(X_{j_k+m}) \right] c_{\xi, j_k, m}^{1-\delta_k},$$

where the sum is taken w.r.t. all  $(\delta_1, \dots, \delta_v)$  with  $\delta_i \in \{0, 1\}$ . Note that all terms in the decomposition (A.30) are Lipschitz. Since  $|c_{\xi, j_k, m}| \leq 4B^2$ , we get

$$\begin{aligned} \|\varphi\|_{\text{Lip}} &\leq \sum_{s=1}^v 2Ls \binom{v}{s} (2B)^{2s-1} (2B)^{2v-2s} = 2Lv(2B)^{2v-1} \sum_{s=1}^v \frac{(v-1)!}{(s-1)!(v-s)!} \\ &= Lv 2^{4v-1} B^{2v-1}. \end{aligned} \quad \blacksquare$$

**Lemma A.5.** Assume (WE)-1. For any  $b, v \geq 1$ ,  $0 = i_1 < \dots < i_b \leq n$  and  $n_k \in \mathbb{N}$ ,  $\sum_{k=1}^b n_k = v$  we define  $g(x) = \mathbb{E}_x[\prod_{k=1}^b f_{i_k}^{n_k}(X_{i_k})]$ , where  $f_{i_k} \in \text{Lip}_{b,d}(L, K)$ . Then  $\|g\|_{\text{Lip}} \leq LvK^{v-1}$ .

*Proof.* Note that  $g(x) = f_0^{n_1}(x)\mathbb{E}_x[\prod_{k=2}^b f_{i_k}^{n_k}(X_{i_k})]$ . We proceed by induction in the number of distinct indices  $b$ . If  $b = 1$ , then, for any  $v \in \mathbb{N}$ ,

$$|f_0^v(x) - f_0^v(y)| = |f_0(x) - f_0(y)| \cdot \left| \sum_{k=0}^{v-1} f_0^k(x)f_0^{v-k-1}(y) \right| \leq vLK^{v-1}d(x, y).$$

Assume  $b > 1$ . Since  $g(x) = f_0^{n_1}(x)P^{i_2}g_1(x)$  with  $g_1(x) = \mathbb{E}_x[\prod_{k=2}^b f_{i_k}^{n_k}(X_{i_k-i_2})]$ ,

$$|g(x) - g(y)| \leq |f_0^{n_1}(x) - f_0^{n_1}(y)| |P^{i_2}g_1(x)| + |f_0^{n_1}(x)| |P^{i_2}g_1(x) - P^{i_2}g_1(y)|.$$

The function  $g_1$  depends on  $b - 1$  indices and  $\sum_{k=2}^b n_k = v - n_1$ . The induction assumption and [18, Theorem 20.1.2] show under (WE)-1 that  $\|P^{i_2}g_1\|_{\text{Lip}} \leq \|g_1\|_{\text{Lip}} \leq L(v - n_1)K^{v-n_1-1}$ . Observe that

$$\|g\|_{\text{Lip}} \leq n_1LK^{n_1-1}K^{v-n_1} + L(v - n_1)K^{v-n_1-1}K^{n_1},$$

and the proof is complete. ■

**A.5. Proof of Proposition 3.7.** We provide the proof only for SGLD, since its adaptation to SGLD-FP is straightforward. Let  $x = (\theta_0^{(1)}, \tilde{S}_0^{(1)})$ ,  $y = (\theta_0^{(2)}, \tilde{S}_0^{(2)})$ . We use the standard synchronous coupling technique adapted from [10, Lemma 1]. Let  $(\xi_k)_{k \geq 0}$  be a sequence of i.i.d.  $d$ -dimensional Gaussian random variables,  $(S_k)_{k \geq 0}$  and  $(\tilde{S}_k)_{k \geq 0}$  be independent mini-batches with  $|S_k| = |\tilde{S}_k| = M$ . Set  $(\theta_0^{(1)}, \theta_0^{(2)}) = (x, y)$ , and define recursively for  $k \geq 0$ ,

$$\begin{aligned} \theta_k^{(i)} &= \theta_{k-1}^{(i)} - \gamma G(\theta_{k-1}^{(i)}, S_k) + \sqrt{2\gamma}\xi_k; \\ G(\theta, S) &= \nabla U_0(\theta) + NM^{-1} \sum_{i \in S} \nabla U_i(\theta). \end{aligned}$$

Finally, define the sequences  $(X_n^{(i)})_{n \geq 0}$ ,  $i = 1, 2$ , as  $X_n^{(i)} = (\theta_n^{(i)}, \tilde{S}_n)$  for any  $n \geq 0$ . Since  $X_k^{(1)}$  and  $X_k^{(2)}$  are distributed according to  $\delta_x \bar{P}^k$  and  $\delta_y \bar{P}^k$ , respectively,

$$W_2^2(\delta_x \bar{P}^k, \delta_y \bar{P}^k) \leq \mathbb{E}[d^2(X_k^{(1)}, X_k^{(2)})] = \mathbb{E}[\|\theta_k^{(1)} - \theta_k^{(2)}\|^2].$$

The rest of the proof follows [10, Lemma 1] and is omitted.

**Acknowledgments.** The authors are greatly indebted to the Associate Editor and the reviewers for the careful reading of the manuscript and pertinent comments. Their constructive feedback helped to improve the quality of this work and shape its final form.

### REFERENCES

- [1] R. ASSARAF AND M. CAFFAREL, *Zero-variance principle for Monte Carlo algorithms*, Phys. Rev. Lett., 83 (1999), pp. 4682–4685.
- [2] J. BAKER, P. FEARNHEAD, E. B. FOX, AND C. NEMETH, *Control variates for stochastic gradient MCMC*, Stat. Comput., 29 (2019), pp. 599–615.
- [3] D. BAKRY AND M. ÉMERY, *Diffusions hypercontractives*, Sémin. Probab. Strasbourg, 19 (1985), pp. 177–206.

- [4] D. BAKRY, I. GENTIL, AND M. LEDOUX, *Analysis and Geometry of Markov Diffusion Operators*, Springer Science & Business Media, New York, 2013.
- [5] D. BELOMESTNY, L. IOSIPOI, E. MOULINES, A. NAUMOV, AND S. SAMSONOV, *Variance reduction for Markov chains with application to MCMC*, *Stat. Comput.*, 30 (2020), pp. 973–997.
- [6] D. BELOMESTNY, L. IOSIPOI, AND N. ZHIVOTOVSKIY, *Variance Reduction via Empirical Variance Minimization: Convergence and Complexity*, preprint, [arXiv:1712.04667](https://arxiv.org/abs/1712.04667) [math. NA], 2017.
- [7] D. V. BELOMESTNY, L. S. IOSIPOI, AND N. K. ZHIVOTOVSKIY, *Variance reduction in Monte Carlo estimators via empirical variance minimization*, *Dokl. Math.*, 98 (2018), pp. 494–497.
- [8] S. BOBKOV AND F. GÖTZE, *Exponential integrability and transportation cost related to logarithmic Sobolev inequalities*, *J. Funct. Anal.*, 163 (1999), pp. 1–28.
- [9] N. BROSE, A. DURMUS, S. MEYN, AND E. MOULINES, *Diffusion Approximations and Control Variates for MCMC*, preprint, [arXiv:1808.01665](https://arxiv.org/abs/1808.01665) [stat. ME], 2018.
- [10] N. BROSE, A. DURMUS, AND E. MOULINES, *The promises and pitfalls of stochastic gradient Langevin dynamics*, in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, 2018, pp. 8278–8288.
- [11] N. S. CHATTERJI, N. FLAMMARION, Y.-A. MA, P. L. BARTLETT, AND M. I. JORDAN, *On the theory of variance reduction for stochastic gradient Monte Carlo*, in *Proceedings of Machine Learning Research*, Vol. 80, 2018, pp. 764–773.
- [12] T. CHEN, E. B. FOX, AND C. GUESTRIN, *Stochastic gradient Hamiltonian Monte Carlo*, in *Proceedings of the 31st International Conference on International Conference on Machine Learning*, ICML’14, 2014, pp. 1683–1691.
- [13] A. S. DALALYAN, *Theoretical guarantees for approximate sampling from smooth and log-concave densities*, *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 79 (2017), pp. 651–676.
- [14] A. S. DALALYAN AND A. G. KARAGULYAN, *User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient*, *Stochastic Process. Appl.*, 129 (2019), pp. 5278–5311.
- [15] A. DEFAZIO, F. BACH, AND S. LACOSTE-JULIEN, *SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives*, in *Advances in Neural Information Processing Systems*, 2014, pp. 1646–1654.
- [16] P. DELLAPORTAS AND I. KONTOYIANNIS, *Control variates for estimation based on reversible Markov chain Monte Carlo samplers*, *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 74 (2012), pp. 133–161.
- [17] H. DJELLOUT, A. GUILLIN, AND L. WU, *Transportation cost-information inequalities and applications to random dynamical systems and diffusions*, *Ann. Probab.*, 32 (2004), pp. 2702–2732.
- [18] R. DOUC, E. MOULINES, P. PRIOURET, AND P. SOULIER, *Markov Chains*, Springer Ser. Oper. Res. Financ. Eng., Springer, Cham, 2018.
- [19] P. DOUKHAN AND S. LOUHICHI, *A new weak dependence condition and applications to moment inequalities*, *Stochastic Process. Appl.*, 84 (1999), pp. 313–342.
- [20] K. A. DUBEY, S. J. REDDI, S. A. WILLIAMSON, B. POZOS, A. J. SMOLA, AND E. P. XING, *Variance reduction in stochastic gradient Langevin dynamics*, in *Advances in Neural Information Processing Systems*, 2016, pp. 1154–1162.
- [21] A. DURMUS AND E. MOULINES, *High-dimensional Bayesian inference via the unadjusted Langevin algorithm*, *Bernoulli*, 25 (2019), pp. 2854–2882.
- [22] J. M. FLEGAL AND G. L. JONES, *Batch means and spectral variance estimators in Markov chain Monte Carlo*, *Ann. Statist.*, 38 (2010), pp. 1034–1070.
- [23] N. FRIEL, A. MIRA, AND C. J. OATES, *Exploiting multi-core architectures for reduced-variance estimation with intractable likelihoods*, *Bayesian Anal.*, 11 (2015), pp. 215–245.
- [24] P. GLASSERMAN, *Monte Carlo Methods in Financial Engineering*, Springer Science & Business Media, New York, 2013.
- [25] E. GOBET, *Monte-Carlo Methods and Stochastic Processes*, CRC Press, Boca Raton, FL, 2016.
- [26] T. E. HANSON, A. J. BRANSCUM, AND W. O. JOHNSON, *Informative g-priors for logistic regression*, *Bayesian Anal.*, 9 (2014), pp. 597–612.
- [27] S. G. HENDERSON, *Variance Reduction via an Approximating Markov Process*, Ph.D. thesis, Stanford University, 1997.
- [28] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance reduction*, in *Advances in Neural Information Processing Systems*, 2013, pp. 315–323.

- [29] M. LEDOUX, *The Concentration of Measure Phenomenon*, Surveys Math. Monogr., American Mathematical Society, Providence, RI, 2001.
- [30] Y.-A. MA, T. CHEN, AND E. FOX, *A complete recipe for stochastic gradient MCMC*, in *Advances in Neural Information Processing Systems*, 2015, pp. 2917–2925.
- [31] K. MARTON, *Bounding  $\bar{d}$ -distance by informational divergence: A method to prove measure concentration*, *Ann. Probab.*, 24 (1996), pp. 857–866.
- [32] S. MENOZZI AND V. LEMAIRE, *On some non asymptotic bounds for the Euler scheme*, *Electron. J. Probab.*, 15 (2010), pp. 1645–1681.
- [33] A. MIRA, R. SOLGI, AND D. IMPARATO, *Zero variance Markov chain Monte Carlo for Bayesian estimators*, *Stat. Comput.*, 23 (2013), pp. 653–662.
- [34] T. NAGAPETRYAN, A. B. DUNCAN, L. HASENCLEVER, S. J. VOLLMER, L. SZPRUCH, AND K. ZYGALAKIS, *The True Cost of Stochastic Gradient Langevin Dynamics*, preprint, [arXiv:1706.02692](https://arxiv.org/abs/1706.02692) [math. NA], 2017.
- [35] R. NICKL AND B. M. PÖTSCHER, *Bracketing metric entropy rates and empirical central limit theorems for function classes of Besov- and Sobolev-type*, *J. Theoret. Probab.*, 20 (2007), pp. 177–199.
- [36] C. J. OATES, M. GIROLAMI, AND N. CHOPIN, *Control functionals for Monte Carlo integration*, *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 79 (2017), pp. 695–718.
- [37] F. OTTO AND C. VILLANI, *Generalization of an inequality by Talagrand and links with the logarithmic Sobolev inequality*, *J. Funct. Anal.*, 173 (2000), pp. 361–400.
- [38] D. J. REZENDE AND S. MOHAMED, *Variational Inference with Normalizing Flows*, preprint, [arXiv:1505.05770](https://arxiv.org/abs/1505.05770) [stat. ML], 2015.
- [39] C. P. ROBERT AND G. CASELLA, *Monte Carlo Statistical Methods*, Springer, New York, 1999.
- [40] G. O. ROBERTS AND R. L. TWEEDIE, *Exponential convergence of Langevin distributions and their discrete approximations*, *Bernoulli*, 2 (1996), pp. 341–363.
- [41] N. L. ROUX, M. SCHMIDT, AND F. R. BACH, *A stochastic gradient method with an exponential convergence rate for finite training sets*, in *Advances in Neural Information Processing Systems*, Vol. 25, 2012, pp. 2663–2671.
- [42] R. Y. RUBINSTEIN AND D. P. KROESE, *Simulation and the Monte Carlo Method*, John Wiley & Sons, Hoboken, NJ, 2016.
- [43] R. SALAKHUTDINOV AND A. MNIH, *Bayesian probabilistic matrix factorization using Markov Chain Monte Carlo*, in *Proceedings of the 25th International Conference on Machine Learning (ICML-08)*, 2008, pp. 880–887.
- [44] M. TALAGRAND, *Transportation cost for Gaussian and other product measures*, *Geom. Funct. Anal.*, 6 (1996), pp. 587–600.
- [45] Y. W. TEH, A. H. THIERY, AND S. J. VOLLMER, *Consistency and fluctuations for stochastic gradient Langevin dynamics*, *J. Mach. Learn. Res.*, 17 (2016), pp. 193–225.
- [46] M. WELLING AND Y. W. TEH, *Bayesian learning via stochastic gradient Langevin dynamics*, in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, 2011, pp. 681–688.
- [47] D. ZOU, P. XU, AND Q. GU, *Subsampled stochastic variance-reduced gradient Langevin dynamics*, in *Proceedings of the International Conference on Uncertainty in Artificial Intelligence*, 2018.