LAPLACIAN SMOOTHING STOCHASTIC GRADIENT MARKOV CHAIN MONTE CARLO*

BAO WANG[†], DIFAN ZOU[‡], QUANQUAN GU[‡], AND STANLEY J. OSHER[§]

Abstract. As an important Markov chain Monte Carlo (MCMC) method, the stochastic gradient Langevin dynamics (SGLD) algorithm has achieved great success in Bayesian learning and posterior sampling. However, SGLD typically suffers from a slow convergence rate due to its large variance caused by the stochastic gradient. In order to alleviate these drawbacks, we leverage the recently developed Laplacian smoothing technique and propose a Laplacian smoothing stochastic gradient Langevin dynamics (LS-SGLD) algorithm. We prove that for sampling from both log-concave and non-log-concave densities, LS-SGLD achieves strictly smaller discretization error in 2-Wasserstein distance, although its mixing rate can be slightly slower. Experiments on both synthetic and real datasets verify our theoretical results and demonstrate the superior performance of LS-SGLD on different machine learning tasks including posterior sampling, Bayesian logistic regression, and training Bayesian convolutional neural networks. The code is available at https://github.com/BaoWangMath/LS-MCMC.

Key words. Langevin dynamics, stochastic gradient, Laplacian smoothing

AMS subject classifications. 60H35, 60J22, 65C05, 65C60

DOI. 10.1137/19M1294356

SIAM J. SCI. COMPUT. Vol. 43, No. 1, pp. A26–A53

1. Introduction. Let $\boldsymbol{x} \in \mathbb{R}^d$ be a machine learning model's parameter with prior distribution $p(\boldsymbol{x})$ and let $p(\boldsymbol{d}|\boldsymbol{x})$ be the likelihood function of the observed data \boldsymbol{d} . Suppose the training data points are generated independently from some unknown distribution; then the posterior distribution of the model parameter \boldsymbol{x} given the entire training dataset $\mathcal{D} = \{\boldsymbol{d}_i\}_{i=1}^n$ is computed as $p(\boldsymbol{x}|\mathcal{D}) \propto p(\boldsymbol{x})\Pi_{i=1}^n p(\boldsymbol{d}_i|\boldsymbol{x})$. While optimization algorithms can be used to find the maximum a posterior point estimator, i.e., $\boldsymbol{x}_{\text{MAP}} = \arg \max_{\boldsymbol{x}} \log p(\boldsymbol{x}|\mathcal{D})$, sampling algorithms such as Langevin dynamics (LD) are used to sample the posterior or the log posterior. In this paper, we consider applying LD-based Markov chain Monte Carlo (MCMC) algorithms to sample $e^{-f(\boldsymbol{x})}$, where

(1.1)
$$f(\boldsymbol{x}) := \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x}) = -\frac{1}{n} \sum_{i=1}^{n} \log p(\boldsymbol{d}_i | \boldsymbol{x}).$$

Here, we normalized the log-likelihood by a factor n for the ease of presentation in the rest of this paper.

^{*}Submitted to the journal's Methods and Algorithms for Scientific Computing section October 21, 2019; accepted for publication (in revised form) October 2, 2020; published electronically January 4, 2021. B. Wang and D. Zou contributed equally; Q. Gu and S. Osher are co-corresponding authors. https://doi.org/10.1137/19M1294356

Funding: This work was supported by the National Science Foundation under grants DMS-1924935, DMS-1952339, DMS-1554564 (STROBE), and SaTC-1717950, by the Air Force Research Laboratory under grants FA9550-18-0167 and MURI FA9550-18-10502, by the Office of Naval Research under grant N00014-18-1-2527, and by the Department of Energy under grant DE-SC0021142.

[†]Scientific Computing and Imaging Institute and Department of Mathematics, University of Utah, Salt Lake City, UT 84112 USA (wangbaonj@gmail.com).

[‡]Department of Computer Science, University of California, Los Angeles, Los Angeles, CA 90095 USA (knowzou@cs.ucla.edu, qgu@cs.ucla.edu).

 $^{^{\$}}$ Department of Mathematics, University of California, Los Angeles, Los Angeles, CA 90095 USA (sjo@math.ucla.edu).

A27

The first-order LD reads

(1.2)

Downloaded 05/07/21 to 155.98.19.70. Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/page/terms

$$d\boldsymbol{X}_t = -\nabla f(\boldsymbol{X}_t) \mathrm{d}t + \sqrt{2\beta^{-1}} \mathrm{d}\boldsymbol{B}_t,$$

where $\mathbf{X}_t \in \mathbb{R}^d$ denotes the point at time t, β denotes the inverse temperature, and $\mathbf{B}_t \in \mathbb{R}^d$ is the standard Brownian term. Under certain assumptions on the negative log posterior (i.e., $f(\mathbf{x})$), the LD (1.2) converges to a unique invariant distribution $\pi \propto e^{-\beta f(\mathbf{x})}$ [14]. Therefore, one can apply an numerical integrator to approximate (1.2) in order to obtain samples that follow the posterior distribution. One simple integrator is to apply the Euler–Maruyama discretization [23] to (1.2), which gives

(1.3)
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \nabla f(\boldsymbol{x}_k) + \sqrt{2\beta^{-1}\eta} \boldsymbol{\epsilon}_k$$

and it is known as the Langevin Monte Carlo (LMC) (a.k.a., unadjusted Langevin algorithm [34]). When the target density, i.e., posterior distribution, is strongly logconcave and log-smooth, Dalalyan and others proved that LMC is able to converge to the target density up to an arbitrarily small sampling error in both total variation and 2-Wasserstein distances [15, 19]. In the same setting, [20] studied the Metropolisadjusted Langevin algorithm, realized by introducing a Metropolis–Hastings (MH) correction step into LMC, and proved its linear convergence rate in total variation distance. Furthermore, the convergence guarantee of LMC for sampling from nonlog-concave distributions has also been established in [35, 40].

Note that the posterior distribution is defined on the whole dataset \mathcal{D} , which is typically extremely large in modern machine learning tasks. Therefore, computing the full gradient $\nabla f(\mathbf{x})$ is inefficient and may dramatically slow down the convergence of sampling algorithms. One solution is to replace the full gradient in (1.3) with a subsampled one, which gives rise to stochastic gradient Langevin dynamics (SGLD) [39]. From the theoretical perspective, the convergence guarantee of SGLD has been proved for both strongly log-concave distributions [16] and non-log-concave distributions [35, 40] in 2-Wasserstein distance. The generalization performance of SGLD for nonconvex optimization has been further investigated in [31]. Although SGLD can drastically reduce the computational cost, it is also observed to have a slow convergence rate due to the large variance caused by the stochastic gradient [36, 37]. In order to reduce the variance of stochastic gradient as well as to improve the convergence rate, Dubey et al. incorporated variance reduction techniques into SGLD [18], which gives rise to a family of variance-reduced LD-based algorithms such as SVRG-LD and SAGA-LD. Chatterji et al. further proved that SVRG-LD and SAGA-LD are able to converge to the target density with fewer stochastic gradient evaluations than SGLD and LMC in certain regimes [8, 43, 44]. However, both SVRG-LD and SAGA-LD require a large amount of extra computation and memory costs, and can only be shown to achieve faster convergence on small to moderate datasets. Therefore, it is natural to ask whether we can reduce the variance of stochastic gradients while maintaining similar computation and memory costs of SGLD.

Recently, Osher et al. integrated Laplacian smoothing and related high-order smoothing techniques into stochastic gradient descent (SGD) to reduce the variance of stochastic gradient on-the-fly [33]. Laplacian smoothing SGD (LSSGD) allows us to take a significantly larger step size than vanilla SGD and reduces the optimality gap in convex optimization when constant step size is used. Empirically, LSSGD preconditions the gradient when the objective function has a large condition number and can avoid local minima. Because of this, LSSGD is applicable to training a large number of deep learning models with good generalization ability. Laplacian smoothing also demonstrates some abilities to avoid saddle point in gradient descent [24]. Most recently, Wang et al. leveraged Laplacian smoothing to improve the utility of machine learning models trained with privacy guarantee [38].

In this paper, we integrate Laplacian smoothing with SGLD, and we call the resulting algorithm Laplacian smoothing SGLD (LS-SGLD). The extra computation of LS-SGLD compared with SGLD is that we need to compute the products of the inverse of two circulant matrices with vectors. We leverage the fast Fourier transform (FFT) to develop fast algorithms to compute these matrix-vector products efficiently, and the resulting algorithms can compute the matrix-vector products with a negligible overhead in both time and memory. Moreover, we prove the convergence rate of LS-SGLD for sampling from both log-concave and non-log-concave densities in 2-Wasserstein distance and show that there exists a trade-off between the discretization error and ergodicity rate. Experimental results show that compared with SGLD, LS-SGLD can achieve much smaller discretization error and similar ergodicity rate, and demonstrate the superior performance of LS-SGLD for a variety of machine learning applications.

1.1. Our contributions. We summarize our main contributions as follows:

- We propose a simple modification on the SGLD, which applies the Laplacian smoothing matrix and its squared root to the stochastic gradient and Gaussian noise vectors, respectively. The continuous and full-gradient counterpart of the modified LS-SGLD has the same stationary distribution as the LD.
- We proposed FFT-based fast algorithms to compute the product of the inverse of circulant matrices with any given vector. By leveraging the structure of eigenvalues and eigenvectors of the circulant matrices, we can compute these products very efficiently with negligible overhead in both time and memory.
- We prove the convergence rate of LS-SGLD for sampling from both logconcave and non-log-concave densities in 2-Wasserstein distance. Specifically, we decompose the sampling error into the discretization error and the ergodicity rate. Moreover, we show that there exists a trade-off between the discretization error and the ergodicity rate of LS-SGLD, as adding Laplacian smoothing can reduce the discretization error but slow down the mixing time.
- We conduct extensive experiments to evaluate the performance of LS-SGLD. First, we show that compared with SGLD, LS-SGLD can achieve a significantly smaller discretization error but similar ergodicity rate, which implies that the overall sampling error of LS-SGLD can be much smaller. Second, we conduct experiments on both synthetic and real data for posterior sampling, Bayesian logistic regression (BLR), and training Bayesian convolutional networks, all of which demonstrate the superior performance of LS-SGLD.

1.2. Additional related work. In addition to the first-order Langevin based algorithms we discussed in the introduction, there also emerges a vast body of work focusing on higher-order Langevin based algorithms. One of the well-known high-order MCMC methods is Hamiltonian Monte Carlo (HMC) [32, 41], which incorporates a Hamiltonian momentum term into the first-order MCMC method in order to improve the mixing time. Other variants of HMC include the reduced-order HMC [1], which has broad applications in scientific computing such as background flow field estimation [6]. Similar to SGLD, a stochastic version of HMC (namely SGHMC) has been further established in [10] and was shown to be able to achieve a faster convergence rate than SGLD in experiments. Ma, Chen, and Fox investigated a family of SGHMC methods and proposed a new state-adaptive sampler on the Riemannian manifold [28]. Chen, Ding, and Carin provided theoretical convergence guarantees of SGHMC in terms of mean square error (MSE) and proposed a second-order symmetric splitting integrator to further improve the discretization error [9]. When the target density is strongly log-concave and log-smooth, Cheng et al. proposed underdamped MCMC and stochastic gradient underdamped MCMC and obtained convergence rates in 2-Wasserstein distance [13]. The convergence rates of these two algorithms have been further established for sampling from non-log-concave densities [12]. Chen et al. proposed the stochastic gradient HMC [11]. However, due to the large variance of stochastic gradients and lack of the MH correction step, SGHMC has also been observed to have a highly biased sampling trajectory [4, 17]. One way to address this issue is to make use of a variance-reduction technique to alleviate the variance of stochastic gradients in SGHMC, which gave rise to stochastic variance-reduced HMC methods [42, 27, 45].

1.3. Organization. We organize this paper as follows. We present LS-SGLD and derive FFT-based fast algorithms for LS-SGLD in section 2. In section 3, we give theoretical guarantees for the performance of LS-SGLD in both log-concave and non-log-concave settings. In section 4, we numerically verify the performance of LS-SGLD on sampling different distributions, training BLR, and convolutional neural nets (CNNs). We conclude this work in section 5.

1.4. Notation. Throughout this paper we use bold uppercase letters \mathbf{A} , \mathbf{B} to denote matrices, bold lowercase letters \mathbf{x} , \mathbf{y} to denote vectors, and lowercase letters \mathbf{x} , y and α , β to denote scalars. We use $x \wedge y$ and $x \vee y$ to denote min $\{x, y\}$ and max $\{x, y\}$, respectively. For continuous-time random vectors, we denote them with italics bold uppercase letters \mathbf{X} , \mathbf{Y} with sub/superscripts. For vector $\mathbf{x} = (x_1, \dots, x_d)^{\top}$, we use $\|\mathbf{x}\|_2 = \sqrt{x_1^2 + \dots + x_d^2}$ to represent its ℓ_2 -norm and use $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^{\top} \mathbf{A} \mathbf{x}}$ to represent its \mathbf{A} -norm, where \mathbf{A} is a semipositive definite matrix. We use $\mathbb{P}(\mathbf{x})$ to denote the distribution of \mathbf{x} , and $\mathcal{W}_2(\cdot, \cdot)$ and $D_{KL}(\cdot \| \cdot)$ denote the 2-Wasserstein distance and Kullback-Leibler (KL) divergence between two distributions, respectively. For a function $f : \mathbb{R}^d \to \mathbb{R}$, we use $\nabla f(\cdot)$ and $\nabla^2 f(\cdot)$ to denote its gradient and Hessian.

2. Algorithms.

2.1. Laplacian smoothing (stochastic) gradient descent. For $\sigma \geq 0$, let $\mathbf{A}_{\sigma} := \mathbf{I} - \sigma \mathbf{L}$ where $\mathbf{I} \in \mathbb{R}^{d \times d}$ and $\mathbf{L} \in \mathbb{R}^{d \times d}$ are the identity and the discrete one-dimensional (1D) Laplacian matrix, respectively. Therefore,

(2.1)
$$\mathbf{A}_{\sigma} := \begin{bmatrix} 1+2\sigma & -\sigma & 0 & \dots & 0 & -\sigma \\ -\sigma & 1+2\sigma & -\sigma & \dots & 0 & 0 \\ 0 & -\sigma & 1+2\sigma & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ -\sigma & 0 & 0 & \dots & -\sigma & 1+2\sigma \end{bmatrix}_{d \times d}$$

To optimize the loss function $f(\mathbf{x}) = 1/n \sum_{i=1}^{n} f_i(\mathbf{x})$, LSSGD [33] takes the following iteration:

(2.2)
$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k - \eta_k \mathbf{A}_{\sigma}^{-1} \nabla f_{i_k}(\boldsymbol{x}^k),$$

where $\eta_k > 0$ is the learning rate, and i_k is a random sample from $[n] := \{1, 2, ..., n\}$. When $\sigma = 0$, LSSGD reduces to SGD. Since \mathbf{A}_{σ} is a circulant matrix, for any vector $\boldsymbol{v}, \mathbf{A}_{\sigma}^{-1}\boldsymbol{v} := \boldsymbol{u}$ can be computed via the FFT in the following way:

$$\mathbf{A}_{\sigma}^{-1}oldsymbol{v} = oldsymbol{u} \Longrightarrow oldsymbol{v} = \mathbf{A}_{\sigma}oldsymbol{u} = oldsymbol{u} - \sigma oldsymbol{d} st oldsymbol{u}$$

where * is the convolution operator, and $\boldsymbol{d} = [-2, 1, 0, \dots, 0, 1]^T$. By the convolution theorem, we have

$$\operatorname{fft}(\boldsymbol{v}) = \operatorname{fft}(\boldsymbol{u}) - \sigma \operatorname{fft}(\boldsymbol{d}) \operatorname{fft}(\boldsymbol{u}).$$

Finally, we arrive at the following FFT-based algorithm for computing $\mathbf{A}_{\sigma}^{-1} \boldsymbol{v}$:

$$\mathbf{A}_{\sigma}^{-1} oldsymbol{v} = \mathrm{ifft}\left(rac{\mathrm{fft}(oldsymbol{v})}{\mathbf{1} - \sigma \cdot \mathrm{fft}(oldsymbol{d})}
ight),$$

where 1 is an all-one vector with the same dimension as v, and the division of two vectors is defined in the coordinatewise way. fft and ifft denote FFT and inverse FFT operators, respectively.

The Laplacian matrix \mathbf{A}_{σ}^{-1} can reduce the variance of stochastic gradient and guarantee at least the same convergence rate as SGD. [33] showed that for an *L*-gradient Lipschitz function $f(\boldsymbol{x})$, i.e., $\|\nabla f(\boldsymbol{x})\|_2 \leq L$, the largest step size for LSSGD is $(1 + 4\sigma)^{1/4}/L$ (with high probability) which is larger than gradient descent's by a factor $(1 + 4\sigma)^{1/4}$.

2.2. Laplacian smoothing Langevin dynamics. We integrate Laplacian smoothing with LD and obtain the following Lapacian smoothing LD (LS-LD):

(2.3)
$$d\boldsymbol{X}_t = -\mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{X}_t) + \sqrt{2\beta^{-1}} \mathbf{A}_{\sigma}^{-1/2} \mathrm{d}\boldsymbol{B}_t.$$

Note that we premultiply the Brownian motion term by $\mathbf{A}_{\sigma}^{-1/2}$ instead of \mathbf{A}_{σ}^{-1} to guarantee that the stationary distribution of the LS-LD remains to be exp $(-\beta f(\boldsymbol{x}))$. We formally state this property in the following proposition.

PROPOSITION 2.1. The stationary distribution, π , of the LS-LD, (2.3), satisfies $\pi \propto e^{-\beta f(\boldsymbol{x})}$.

The proof of Proposition 2.1 can be found in Appendix A. If we apply the Euler-Maruyama scheme to discretize (2.3), we end up with the following discrete algorithm, namely Laplacian smoothing gradient Langevin dynamics (LS-GLD):

(2.4)
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}_k) + \sqrt{2\beta^{-1}\eta} \mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}_k,$$

where $\epsilon_k \sim N(\mathbf{0}, \mathbf{I}_{d \times d})$. In practice, we use the minibatch gradient $\mathbf{g}_k = \sum_{i \in \mathcal{I}_k} \nabla f_i(\mathbf{x}_k) / |\mathcal{I}_k|$ with $\mathcal{I}_k \subset [n]$ to replace the gradient in (2.4), and we arrive at the following LS-SGLD:

(2.5)
$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \mathbf{A}_{\sigma}^{-1} \mathbf{g}_k + \sqrt{2\beta^{-1}\eta} \mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}_k.$$

We summarize LS-SGLD in Algorithm 2.1. In the remaining part of this section, we will present FFT-based fast algorithms for computing $\mathbf{A}_{\sigma}^{-1}\mathbf{g}_{k}$ and $\mathbf{A}_{\sigma}^{-1/2}\boldsymbol{\epsilon}_{k}$.

2.3. FFT-based implementation of LS-SGLD.

2.3.1. Circulant matrix and convolutional operation. In this subsection, we list a few results on the circulant matrix which will be the basic recipes for designing an FFT-based algorithm for solving (2.4).

A30

 $\begin{array}{l} \label{eq:algorithm 2.1. LS-SGLD.} \\ \hline \textbf{Input: Training data, learning rate η, minibatch size B, inverse temperature β, Laplacian smoothing constant σ. \\ \hline \textbf{Initialization: Set $x_0 = 0$.} \\ \textbf{for $k = 0, 1, \ldots, K-1$ do} \\ \textbf{Uniformly sample $\mathcal{I}_k \subset [n]$ with $|\mathcal{I}_k| = B$.} \\ \textbf{Compute the minibatch stochastic gradient $\sum_{i \in \mathcal{I}_k} \nabla f_i(\pmb{x}_k)/B$.} \\ \pmb{x}_{k+1} = \pmb{x}_k - \eta \mathbf{A}_{\sigma}^{-1} \mathbf{g}_k + \sqrt{2\beta^{-1}\eta} \mathbf{A}_{\sigma}^{-1/2} \pmb{\epsilon}_k, \text{ where $\mathbf{\epsilon}_k \sim N(\mathbf{0}, \mathbf{I}_{d \times d})$.} \\ \textbf{Output: x_0, \ldots, x_K.} \end{array}$

LEMMA 2.2 (see [21]). The normalized eigenvectors of the $d \times d$ circulant matrix,

(2.6)
$$\mathbf{C} = \begin{bmatrix} c_0 & c_{d-1} & \dots & c_2 & c_1 \\ c_1 & c_0 & c_{d-1} & \dots & c_2 \\ \dots & \dots & \dots & \dots & \dots \\ c_{d-2} & \dots & \dots & \dots & c_{d-1} \\ c_{d-1} & c_{d-2} & \dots & c_1 & c_0 \end{bmatrix},$$

are given by

$$\mathbf{v}_j = \frac{1}{\sqrt{d}} \left(1, w_j, w_j^2, \dots, w_j^{n-1} \right), \quad j = 0, 1, \dots, d-1,$$

where $w_j = \exp(i\frac{2\pi j}{d})$ are the *j*th roots of unity and *i* is the imaginary unit. The corresponding eigenvalues are then given by

$$\lambda_j = c_0 + c_{d-1}w_j + c_{d-2}w_j^2 + \dots + c_1w_j^{d-1}, \quad j = 0, 1, \dots, d-1$$

LEMMA 2.3 (see [21]). The inverse of a circulant matrix is circulant.

LEMMA 2.4 (see [21]). The square root of a circulant matrix is circulant.

LEMMA 2.5. For any circulant matrix **C** of the form in (2.6), and for any given vector \mathbf{v} , let $\mathbf{u} = \mathbf{C}^{-1}\mathbf{v}$; then \mathbf{u} can be computed by the FFT with sublinear scaling in the following way:

(2.7)
$$\boldsymbol{u} = \operatorname{ifft}\left(\frac{\operatorname{fft}(\boldsymbol{v})}{\operatorname{fft}(\boldsymbol{c})}\right),$$

where \mathbf{c} is the first row of the matrix \mathbf{C} , and the division in (2.7) is defined coordinatewise.

Proof. Since **C** is a circulant matrix we have $v = \mathbf{C}u = c * u$; therefore $\mathrm{fft}(v) = \mathrm{fft}(c) \cdot \mathrm{fft}(u)$.

2.3.2. Fast algorithm for computing the square root of Laplacian smoothing. We will derive an FFT-based algorithm for computing $\mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}_{k}$ in this subsection. According to Lemmas 2.3 and 2.4, $\mathbf{A}_{\sigma}^{-1/2}$ is circulant. Note \mathbf{A}_{σ}^{-1} is positive definite; we denote its eigen-decomposition as

$$\mathbf{A}_{\sigma}^{-1} = \mathbf{Q} \Lambda \mathbf{Q}^{-1},$$

where $\mathbf{Q} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d]^T$ with \mathbf{v}_i being the eigenvector associated with the eigenvalue $\lambda_i > 0$, and $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$. Therefore, we have

(2.8)
$$\mathbf{A}_{\sigma}^{-1/2} = \mathbf{Q}\sqrt{\Lambda}\mathbf{Q}^{-1}$$

where $\sqrt{\Lambda} = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_d}).$

Furthermore, note that \mathbf{A}_{σ} is symmetric, therefore $\mathbf{Q}^{-1} = \mathbf{Q}^{T}$. It follows that we can compute $\mathbf{A}_{\sigma}^{-1/2}$ without inverting the matrix \mathbf{Q} . By the fact that $\mathbf{A}_{\sigma}^{-1/2}$ is circulant, we have $\mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}_{k} = \operatorname{ifft}(\operatorname{fft}(\boldsymbol{b}) \cdot \operatorname{fft}(\boldsymbol{\epsilon}_{k}))$, where \boldsymbol{b} is the first row of $\mathbf{A}_{\sigma}^{-1/2}$.

Remark 2.6. In computing (2.8), there is no need to store the matrix \mathbf{Q} ; according to Lemma 2.2, each row of \mathbf{Q} and $\sqrt{\Lambda}$ can be written explicitly, which enables us to compute $\mathbf{A}_{\sigma}^{-1/2}$ quickly with negligible memory overhead and scalable to very high dimensional problems.

3. Main theoretical results. We first make the following three assumptions regarding the function f(x).

Assumption 1 (dissipativeness). For any $x \in \mathbb{R}^d$, there exist constants m and b such that

$$\langle \nabla f(\boldsymbol{x}), \boldsymbol{x} \rangle \ge m \|\boldsymbol{x}\|_2^2 - b.$$

This assumption has been widely made to study the convergence of Langevin-based sampling algorithms [30, 35, 40, 44], which is essential to guarantee the convergence of the continuous-time Langevin dynamics (1.2).

Assumption 2 (smoothness). For any $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$, there exists a positive constant M such that for all i = 1, ..., n, it holds that

$$\|\nabla f_i(\boldsymbol{x}) - \nabla f_i(\boldsymbol{y})\|_2 \leq M \|\boldsymbol{x} - \boldsymbol{y}\|_2.$$

Unlike Assumption 1, Assumption 2 is made for all component functions $f_i(x)$.

Assumption 3 (bounded variance). For any $\boldsymbol{x} \in \mathbb{R}^d$, there exists a constant ω such that the variance of the stochastic gradient is bounded as follows:

$$\mathbb{E}[\|\nabla f_i(\boldsymbol{x}) - \nabla f(\boldsymbol{x})\|_2^2] \le d\omega^2$$

DEFINITION 3.1 (logarithmic Sobolev inequality). Let μ be a probability measure; then we say μ satisfies the logarithmic Sobolev inequality with constant λ if for any smooth function g, the following holds:

$$\int g^2 \log g^2 d\mu - \int g^2 d\mu \log \int g^2 d\mu \le \lambda \int \|\nabla g\|_2^2 d\mu.$$

Then the following proposition states that if the function $f(\cdot)$ satisfies Assumptions 1 and 2, the target density $\pi \propto e^{-f(\boldsymbol{x})}$ satisfies the logarithmic Sobolev inequality.

PROPOSITION 3.2 (see [35]). Under Assumptions 1 and 2, the target density $\pi \propto e^{-f(\boldsymbol{x})}$ satisfies the logarithmic Sobolev inequality with some constant $\lambda > 0$.

It has been shown in [19, 35] that if the function $f(\mathbf{x})$ is smooth and strongly convex (which is stronger than Assumption 1), the logarithmic Sobolev constant λ is a universal constant. However, if the function $f(\mathbf{x})$ is nonconvex, in the worst case the logarithmic Sobolev constant λ has exponential dependency on the problem dimension d and inverse temperature β [7, 35].

A32

3.1. Convergence analysis of sampling from log-concave densities. In this subsection, we assume that the target density is log-concave, which is equivalent to the following assumption on the function $f(\mathbf{x})$.

Assumption 4 (convexity). For any $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$, it holds that

$$f(\boldsymbol{x}) - f(\boldsymbol{y}) \ge \langle \nabla f(\boldsymbol{y}), \boldsymbol{x} - \boldsymbol{y} \rangle.$$

Then we are ready to establish the convergence rate of LS-SGLD for sampling from log-concave densities, which is stated in the following theorem.

THEOREM 3.3. Under Assumptions 1, 2, 3, and 4, if we set the step size $\eta \leq Cm\beta^{-1}/M^2$ for some sufficiently small constant C, there exist constants $c_0 \in [||\mathbf{A}_{\sigma}||_2^{-1}]$, 1], $\gamma_1 \in [||\mathbf{A}_{\sigma}||_2^{-2}, 1]$, and $\gamma_2 = d^{-1} \sum_{i=1}^d (1+2\sigma-2\sigma\cos(2\pi i/d))^{-1}$ such that the output of LS-SGLD satisfies

(3.1)
$$\mathcal{W}_{2}(\mathbb{P}(\boldsymbol{x}_{K}), \pi) \leq \left(\frac{2\gamma_{1}K\eta^{2}\beta d\omega^{2}}{B}\right)^{1/2} + \left[8\gamma_{2}K\eta^{2}(K+1)\beta d\eta\right]^{1/2} + \left[2\lambda\left(\beta f(0) + \log(\Lambda)\right)\right]^{1/2}e^{-c_{0}K\eta/(2\beta\lambda)},$$

where $\Lambda = \int_{\mathbb{R}^d} e^{-\beta f(\boldsymbol{x})} d\boldsymbol{x}$ and λ denotes the logarithmic Sobolev constant of the target distribution $\pi \propto e^{-\beta f(\boldsymbol{x})}$.

Remark 3.4. We emphasize that the three terms on the right-hand side (R.H.S.) of (3.1) have their respective meanings. In particular, the first and second terms represent the discretization errors introduced by the stochastic gradient estimator and the numerical integrator of (2.3), respectively. The third term represents the ergodicity of the continuous-time Markov process (2.3), which characterizes the mixing time of LS-LD (2.3). Moreover, we remark here that the convergence rate of LS-GLD (LS-SGLD with full gradient) can be directly implied from Theorem 3.3 by removing the first term on the R.H.S. of (3.1).

Based on Theorem 3.3, we can also derive the convergence rate of SGLD in the same setting by setting $\mathbf{A}_{\sigma} = \mathbf{I}$ (i.e., $\sigma = 0$), which implies that the constants γ_1, γ_2 , and c_0 in Theorem 3.3 are all 1's. We formally state the convergence result of SGLD in the following corollary.

COROLLARY 3.5. Under the same assumptions in Theorem 3.3, the output of standard SGLD, denoted by \mathbf{y}_K , satisfies

(3.2)
$$\mathcal{W}_{2}(\mathbb{P}(\boldsymbol{y}_{K}), \pi) \leq \left(\frac{2K\eta^{2}d\omega^{2}}{B}\right)^{1/2} + \left[8K\eta^{2}(K+1)\beta^{-1}d\eta\right]^{1/2} + \left[2\lambda\left(\beta f(0) + \log(\Lambda)\right)\right]^{1/2}e^{-K\eta/(2\beta\lambda)}.$$

Remark 3.6. We can now compare the convergence rates of LS-SGLD and SGLD. In terms of the discretization error, it is clear that LS-SGLD is strictly better since the constants γ_1 and γ_2 are strictly less than 1 (some values of γ_2 corresponding to different choices of σ and d can be found in Table 1). In terms of the ergodicity of the continuous-time Markov process (the third terms in (3.1) and (3.2)), LS-SGLD is worse than SGLD due to the fact that $c_0 \leq 1$. Therefore, there exists a trade-off between the discretization error and the ergodicity rate of LS-SGLD. In our experiments we will conduct numerical evaluations of these error terms and demonstrate that LS-LD and LD achieve similar ergodicity performance (i.e., mixing time), but LS-SGLD can achieve a significantly smaller discretization error. A34

(3.3)

TABLE 1						
The values	of γ_2	corresponding	to	some	σ	$and \ d.$

σ	1	2	3	4	5
d = 1000 d = 10000 d = 100000	$0.268 \\ 0.268 \\ 0.268$	$\begin{array}{c} 0.185 \\ 0.185 \\ 0.185 \end{array}$	$0.149 \\ 0.149 \\ 0.149$	$0.128 \\ 0.128 \\ 0.128$	$\begin{array}{c} 0.114 \\ 0.114 \\ 0.114 \end{array}$

3.2. Convergence analysis of sampling from non-log-concave densities. Here we consider the setting where the target density is no longer log-concave. The following theorem states the convergence rate of LS-SGLD in 2-Wasserstein distance.

THEOREM 3.7. Under Assumptions 1, 2, and 3, if set the step size $\eta \leq Cm\beta^{-1}/M^2$ for some sufficiently small constant C, there exist constants $c_0 \in [\|\mathbf{A}_{\sigma}\|_2^{-1}, 1], \gamma_1 \in [\|\mathbf{A}_{\sigma}\|_2^{-2}, 1], \gamma_2 = d^{-1} \sum_{i=1}^d (1+2\sigma-2\sigma\cos(2\pi i/d))^{-1}, \text{ and } \bar{\Gamma} = (3/2+2(b+\beta^{-1}d))^{1/2}$ such that the output of LS-SGLD satisfies

$$\mathcal{W}_{2}\left(\mathbb{P}(\boldsymbol{x}_{K}),\pi\right) \leq \bar{\Gamma}(K\eta)^{1/2} \left[\left(\frac{\gamma_{1}\beta d\omega^{2}}{B}K\eta + 2\gamma_{2}M^{2}dK\eta^{2}\right)^{1/2} + \left(\frac{\gamma_{1}\beta d\omega^{2}}{B}K\eta + 2\gamma_{2}M^{2}dK\eta^{2}\right)^{1/4} \right] + \left[2\lambda\left(\beta f(0) + \log(\Lambda)\right)\right]^{1/2}e^{-c_{0}K\eta/(2\beta\lambda)},$$

where $\Lambda = \int_{\mathbb{R}^d} e^{-\beta f(\boldsymbol{x})} d\boldsymbol{x}$ and λ denotes the logarithmic Sobolev constant of the target distribution $\pi \propto e^{-\beta f(\boldsymbol{x})}$.

Remark 3.8. The convergence rate of SGLD in 2-Wasserstein distance can also be obtained from Theorem 3.7 by setting $\mathbf{A}_{\sigma} = \mathbf{I}$, which implies that the constants c_0, γ_1, γ_2 become all 1's. It can be verified that the resulting convergence rate matches that proved in [35]. As a clear comparison, the discretization error induced by both the stochastic gradient and the numerical integrator of LS-SGLD (the first bracket term of (3.3)) is smaller than that of SGLD, while the ergodicity term of LS-SGLD (the last term of (3.3)) is worse than that of SGLD. Again, we will experimentally demonstrate that the mixing time of LS-LD is not much slower compared with LD, but LS-SGLD can achieve significantly smaller discretization error than SGLD.

4. Numerical results. In this section, we will perform numerical experiments on sampling 2D distributions, training BLR, and training CNNs. Throughout all the experiments, we regard SGLD [39] and preconditioned SGLD (pSGLD) [26], which considers local curvature of f(x) with RMSProp type adaptive step size, as benchmarks. In addition, we also incorporated the precondition technique, proposed in [26], into LS-SGLD, which leads to a variant of LS-SGLD, namely Laplacian smoothing preconditioned SGLD (LS-pSGLD).

4.1. Numerical simulations on synthetic datasets.

4.1.1. 2D Gaussian distribution. As a simple illustration, we apply the proposed LS-SGLD and LS-pSGLD to sample a 2D Gaussian distribution, studied in [10], with the probability density function $e^{f(\boldsymbol{x})} = \exp(\frac{1}{2}\boldsymbol{x}^T\Sigma^{-1}\boldsymbol{x})$ with $\boldsymbol{x} \in \mathbb{R}^2$ where $\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$. We let the prior be $p(\boldsymbol{x}) = \mathcal{N}(0, \nu^2 \mathbf{I})$ with $\nu = 1$. For both SGLD and pSGLD, we let the step size be 0.19, which is obtained based on the grid search. For LS-SGLD and LS-pSGLD, we let the Laplacian smoothing parameter σ be 0.1 with step size either 0.19 or $0.19(1 + 4\sigma)^{1/4}$. It is worth noting that in 2D \mathbf{A}_{σ} becomes

LAPLACIAN SMOOTHING STOCHASTIC GRADIENT MCMC

(4.1)
$$\mathbf{A}_{\sigma} = \begin{bmatrix} 1+\sigma & -\sigma \\ -\sigma & 1+\sigma \end{bmatrix}$$

To measure the quality of samples, we consider the MSE between the true and reconstructed covariance matrices, and we calculate the autocorrelation time of the samples to verify the efficacy of different samplers in sampling the correlated distribution above. The autocorrelation time is defined as

(4.2)
$$\tau = \frac{1}{2} + \sum_{t=1}^{\infty} \frac{A(t)}{A(0)},$$

where $A(t) = \mathbb{E}[(\bar{\phi}_{\eta} - \phi(\boldsymbol{x}_{0}))(\bar{\phi}_{\eta} - \phi(\boldsymbol{x}_{t}))]$ for any given bounded function $\phi(\boldsymbol{x})$, $\bar{\phi} = \int_{\chi} \phi(\boldsymbol{x}) p(\boldsymbol{x}|\mathcal{D}) d\boldsymbol{x}$ is the population mean of ϕ , and the empirical mean $\hat{\phi} = \frac{1}{S_{T}} \sum_{t=1}^{T} \eta_{t} \phi(\boldsymbol{x}_{t})$ with η_{t} being the step size at the *t*th step and $S_{T} = \sum_{t=1}^{T} \eta_{t}$.

Figures 1(a) and (c) plot the first 600 samples from the target distribution by different samplers. We use the same step size 0.19 for all four samplers in the experiments shown in Figure 1(a) and use a larger step size $0.19(1 + 4\sigma)^{1/4}$ for LS-SGLD and LS-pSGLD in experiments shown in Figure 1(c). Qualitatively, Laplacian smoothing can enhance the quality of samples, and the improvement becomes more remarkable when we use a larger step size. Next, we draw 2×10^5 samples from the target distribution by different samplers and we use these samples to reconstruct the covariance matrix. For SGLD and pSGLD, we use a set of step size



FIG. 1. Contrasting sampling of a 2D Gaussian distribution with covariance matrix Σ , with $\Sigma_{11} = \Sigma_{22} = 1$ and $\Sigma_{12} = \Sigma_{21} = 0.9$, using different samplers. (a) and (c) The first 600 samples draw by SGLD, pSGLD. (b) and (d) LS-SGLD, and LS-pSGLD. In (c) and (d), we multiply the step size for LS-SGLD and LS-pSGLD by a factor $(1 + 4\sigma)^{1/4}$.

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

{0.19, 0.19 × 0.8, 0.19 × 0.8², 0.19 × 0.8³, 0.19 × 0.8⁴}. For LS-SGLD and LS-pSGLD we test two sets of step sizes: (i) {0.19, 0.19 × 0.8, 0.19 × 0.8², 0.19 × 0.8³, 0.19 × 0.8⁴}, (ii) {0.19, 0.19 × 0.8, 0.19 × 0.8², 0.19 × 0.8³, 0.19 × 0.8⁴} × (1 + 4 σ)^{1/4}. Figures 1(b) and (d) plot the autocorrelation time versus reconstruction error of the covariance matrix. In (b) we use the same set of step sizes for all four samplers, and in (d) we use a larger step size for LS-SGLD and LS-pSGLD. We see that reconstruction errors can be reduced significantly when Laplacian smoothing is used. Moreover, Laplacian smoothing can also reduce the autocorrelation time in pSGLD.

4.1.2. 2D Gaussian mixture distribution. In this subsection, we compare the performance of SGLD, pSGLD, LS-SGLD, and LS-pSGLD on a Gaussian mixture distribution. In particular, we consider the target distribution

$$\pi \propto \exp\left(-f(\boldsymbol{x})\right) = \exp\left(-\frac{1}{n}\sum_{i=1}^{n}f_{i}(\boldsymbol{x})\right), \quad n = 500,$$

where each component $\exp\left(-f_i(\boldsymbol{x})\right)$ is defined as

$$\exp\left(-f_i(\boldsymbol{x})\right) = \frac{2}{3}e^{-\frac{\|\boldsymbol{x}-\boldsymbol{a}_i\|_2^2}{2}} + \frac{1}{3}e^{-\frac{\|\boldsymbol{x}+\boldsymbol{a}_i\|_2^2}{2}},$$

where we sample a_i by the MCMC sampler with MH correction from the following 2D Gaussian distribution:

$$\mathcal{N}\left(\begin{bmatrix}2\\2\end{bmatrix},\begin{bmatrix}2&0\\0&2\end{bmatrix}\right).$$

The function $f_i(\mathbf{x})$ and its gradient can be simplified as

$$egin{aligned} f_i(oldsymbol{x}) &= rac{\|oldsymbol{x} - oldsymbol{a}_i\|_2^2}{2} - \log\left(rac{2}{3} + rac{1}{3}\exp\left(-2\langleoldsymbol{a}_i,oldsymbol{x}
angle
ight)
ight), \
onumber \nabla f_i(oldsymbol{x}) &= oldsymbol{x} - oldsymbol{a}_i + rac{2oldsymbol{a}_i}{2 + \exp\left(2\langleoldsymbol{x},oldsymbol{a}_i
angle
ight)}. \end{aligned}$$

It can be easily verified that if a_i satisfies $||a_i||_2 > 3/2$, the function $f_i(x)$ defined above is nonconvex. Moreover, it can be seen that

$$\langle
abla f_i(oldsymbol{x}),oldsymbol{x}
angle = \|oldsymbol{x}\|_2^2 - rac{\exp\left(2\langleoldsymbol{x},oldsymbol{a}_i
angle
ight)}{2+\exp\left(2\langleoldsymbol{x},oldsymbol{a}_i
angle
ight)}\langleoldsymbol{a}_i,oldsymbol{x}
angle \geq rac{1}{2}\|oldsymbol{x}\|_2^2 - rac{1}{2}\|oldsymbol{a}_i\|_2^2,$$

which suggests that the function $f_i(\boldsymbol{x})$ satisfies the dissipative Assumption 1 with $m = \frac{1}{2}$ and $b = \frac{\|\boldsymbol{a}_i\|_2^2}{2}$, and it further implies that $f(\boldsymbol{x})$ is also dissipative.

Since it takes a large number of samples to characterize the distribution, which makes repeated experiments computationally expensive, we instead follow [3, 44, 45] to use iterates along one Markov chain to visualize the distribution of iterates obtained by MCMC algorithms. We run the four samplers with different numbers of iterations where we set the batch size to be 10. We plot the distributions generated by different samplers with different numbers of iterations in Figure 2. As shown in Figures 2(c), (f), and (i), when the number of iterations is large enough, e.g., 10^6 , the sample distributions of all three samplers matches well with the reference distribution (sampled by ground-truth sampler, e.g., MCMC with MH step). However, when the number of iterations is not enough, there is a large discrepancy between the sample and target distributions, as shown in Figures 2(a), (d), and (g). With a moderate

LAPLACIAN SMOOTHING STOCHASTIC GRADIENT MCMC



FIG. 2. Kernel density plots of samples generated from Gaussian mixture distribution using SGLD, LS-SGLD, pSGLD, and LS-pSGLD. We set $\sigma = 1.0$ for LS-SGLD and LS-pSGLD.

Table 2

2-Wasserstein distance between samples sampled by MCMC with MH correction and different SGLD.

# of samples	1E5	$5\mathrm{E5}$	9E5
${ m SGLD} m _pSGLD m _LS-SGLD$	$0.695 \\ 5.364 \\ 0.421$	$6.726 \\ 0.286 \\ 0.414$	$0.285 \\ 6.728 \\ 0.418$

number of iterations, say, 5×10^5 , the sample distribution from LS-SGLD is better than the other two (Figures 2(b), (e), and (h)).

Let us further evaluate the sample quality in a quantitative approach. We first apply the MCMC with MH step to sample 10K samples from the above target distribution. Then we apply SGLD, pSGLD, and LS-SGLD to sample different numbers of samples, respectively, from the target distribution. We measure the 2-Wasserstein distance between the last 10K samples of the different number of samples by the above three stochastic gradient samplers with the MH samples. We list the Wasserstein distance between the last 10K samples of different numbers of samples from different samplers with the MH samples in Table 2. These results show that the samples generated by LS-SGLD are consistently closer to the samples from MCMC with MH correction.

4.1.3. Comparison of the mixing time between LD and LS-LD. To verify that Laplacian smoothing in practice does not slow down the mixing rate of the continuous-time Markov process, we conduct the following experiments. First, we apply the MCMC with MH correction step to sample 10K points, respectively, from the following two distributions:



FIG. 3. MSE between the true and reconstructed means from different numbers of samples generated by LD and LS-LD (10 independent runs).

• Gaussian distribution with the probability density function

4.3)
$$p(x,y) = \frac{1}{9\pi} \exp\left(-\left(\frac{(x-1)^2}{3^2} + \frac{(y-2)^2}{3^2}\right)\right),$$

• the Gaussian mixture distribution described in subsection 4.1.2.

Second, we use either LD or LS-LD (which can be approximated by Euler–Maruyama discretization with very small step size), to draw samples from the above two distributions and use these samples to estimate the mean of the target densities. Figure 3 plots the MSE between the true and reconstructed (from a different number of samples) means, and they show that LD and LS-LD perform similarly in reconstructing the mean of the target densities.

4.2. Bayesian logistic regression. Suppose we observe n independent and identically distributed samples $\{d_i, y_i\}_{i=1,2,...,n}$ where $d_i \in \mathbb{R}^d$ and $y_i \in \{-1, 1\}$ denote the feature and the corresponding label of the *i*th sample instance. The likelihood of the BLR model is given by

$$p(y_i|\boldsymbol{d}_i, \boldsymbol{x}) = \frac{1}{1 + \exp\left(-y_i\langle \boldsymbol{d}_i, \boldsymbol{x}\rangle\right)}$$

where \boldsymbol{x} is the parameter to be learned. We use a Gamma prior $p(\boldsymbol{x}) \propto \|\boldsymbol{x}\|_2^{-\lambda} \exp(-\theta \|\boldsymbol{x}\|_2)$ with $\lambda = 1$ and $\theta = 10^{-2}$. Then we formulate the logarithmic posterior distribution as follows:

$$\log\left[p(oldsymbol{x}|oldsymbol{d}_1,oldsymbol{d}_2,\ldots,oldsymbol{d}_n;y_1,y_2,\ldots,y_n)
ight] \propto -rac{1}{n}\sum_{i=1}^n f_i(oldsymbol{x}),$$

where $f_i(\boldsymbol{x}) = n \log \left(1 + e^{-y_i \langle \boldsymbol{d}_i, \boldsymbol{x} \rangle}\right) + \lambda \log \left(\|\boldsymbol{x}\|_2\right) + \theta \|\boldsymbol{x}\|_2.$

We use SGLD, pSGLD, LS-SGLD, and LS-pSGLD with batch size 5 to train a BLR model on the benchmark a3a dataset from the UCI machine learning repository.¹ The a3a dataset contains 3185 training data and 29376 test instances; each data instance is of dimension 122. We use the grid search to determine the optimal learning rate for SGLD (0.001) and pSGLD (0.002), and then we multiply them by $(1 + 4\sigma)^{1/4}$ to get the learning rate for LS-SGLD and LS-pSGLD. We set the burn-in to be 1000 for all four samplers. After burn-in, we compute the moving average of the sample parameters to estimate the regression parameters \boldsymbol{x} . We plot iteration versus negative

(

¹https://archive.ics.uci.edu/ml/index.php.



FIG. 4. Convergence comparison for BLR, where the X-axis represents the number of iterations and the Y-axis represents the negative log-likelihood/accuracy. (a) Negative log-likelihood on training dataset; (b) accuracy on training dataset; (c) negative log-likelihood on test dataset; (d) accuracy on test dataset.

log-likelihood and accuracy in Figure 4, and we see that Laplacian smoothing reduces the negative log-likelihood and increases the accuracy. The preconditioning accelerates mixing initially; however, the gap between sampling and target distribution is remarkably larger than the case without preconditioning.

4.2.1. Variance reduction in stochastic gradient. We numerically verify the efficiency of variance reduction on BLR for a3a dataset classification. We first compute a path by full batch SGLD with the same learning rate as before, and meanwhile, we record the Laplacian smoothing gradient on each point along the path. Then we compute the Laplacian smoothing stochastic gradients on each point along the path by using different batch size and σ . We run 100 independent experiments to acquire the Laplacian smoothing stochastic gradients, and then we compute the variance of these stochastic gradients by regarding the full batch Laplacian smoothing gradient as the mean. In Table 3, we report the maximum variance, among all coordinates of the gradient and all points on the descent path, for each pair of batch size and σ .

TABLE 3 The maximum variance of the stochastic gradients generated by LS-SGLD on training BLR on the a3a data. $\sigma = 0$ reduces to SGLD.

Batch size	10	15	50
$\sigma = 0$ $\sigma = 0.5$ $\sigma = 1.0$ $\sigma = 2.0$	7.69E-1 2.56E-1 1.54E-1 8.52E-2	3.17E-1 1.06E-2 6.37E-2 3.54E-2	5.69E-2 1.96E-2 1.21E-2 7.04E-3



FIG. 5. Convergence comparison for Bayesian convolutional neural network, where the X-axis represents the number of iterations and the Y-axis represents the negative log-likelihood/accuracy. (a) Negative log-likelihood on training dataset; (b) accuracy on training dataset; (c) negative log-likelihood on test dataset; (d) accuracy on test dataset.

4.3. Bayesian convolutional neural network. We consider training a CNN by SGLD, pSGLD, LS-SGLD, and LS-pSGLD on the MNIST benchmark with batch size 100; the architecture of the CNN is

CNN: $\operatorname{input}_{28 \times 28} \to \operatorname{conv}_{20,5,2} \to \operatorname{conv}_{20,20,5} \to \operatorname{fc}_{128} \to \operatorname{softmax}$.

The notation $\operatorname{conv}_{c,k,m}$ denotes a 2D convolutional layer with c output channels, each of which is the sum of a channelwise convolution operation on the input using a learnable kernel of size $k \times k$; it further adds ReLU nonlinearity and max-pooling with stride size m. fc₁₂₈ is an affine transformation that transforms the input to a vector of dimension 128. Finally, the tensors are activated by a multiclass logistic function.

Similar to BLR, we use a Gamma prior $p(\boldsymbol{x}) \propto \|\boldsymbol{x}\|_2^{-\lambda} \exp(-\theta \|\boldsymbol{x}\|_2)$ with $\lambda = 1$ and $\theta = 5e^{-4}$. Again, we use the grid search to find the optimal step size for SGLD and pSGLD which is 0.02 and 2e - 4, respectively. We multiply the optimal step size for SGLD and pSGLD by a factor $(1 + 4\sigma)^{1/4}$ to get the step size for LS-SGLD and LS-pSGLD, and we let $\sigma = 0.5$ for Laplacian smoothing. The comparisons between different sampling algorithms are plotted in Figure 5, where we see that Laplacian smoothing reduces the negative log-likelihood and increases the accuracy of both training and test datasets. The preconditioning accelerates mixing and reduces the gap between sampling and target distribution. Here, we applied early stopping in training CNN by pSGLD and LS-pSGLD.

5. Conclusions. In this paper, we integrate Laplacian smoothing with SGLD to reduce the gap between the sample and target distributions. The resulting algorithm also allows us to take a larger step size. The proposed algorithm is simple to implement and the extra computation and memory costs compared with the SGLD are negligible

A41

when the FFT-based algorithms are employed to resolve the dynamics of the resulting LS-SGLD. We show, both theoretically and empirically, that LS-SGLD can improve the sample quality. It is straightforward to extend Laplacian smoothing to the other MCMC algorithms, e.g., the stochastic gradient HMC [10].

Appendix A. Missing proof in section 2. In this section, we provide the proof of Proposition 2.1.

Proof of Proposition 2.1. Let $p(\boldsymbol{x},t)$ be the distribution of \boldsymbol{X}_t . Then we know that $p(\boldsymbol{x},t)$ satisfies the following Fokker-Planck equation:

(A.1)
$$\begin{aligned} \frac{\partial p(\boldsymbol{x},t)}{\partial t} &= \frac{1}{\beta} \langle \mathbf{A}_{\sigma}^{-1}, \nabla^2 p(\boldsymbol{x},t) \rangle + \langle \nabla, p(\boldsymbol{x},t) \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}) \rangle \\ &= \frac{1}{\beta} \langle \nabla, \mathbf{A}_{\sigma}^{-1} \nabla p(\boldsymbol{x},t) \rangle + \langle \nabla, p(\boldsymbol{x},t) \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}) \rangle, \end{aligned}$$

where $\langle \nabla, \mathbf{h}(\boldsymbol{x}) \rangle$ denotes the divergence of the vector field $\mathbf{h}(\boldsymbol{x})$. Since the stationary distribution π satisfies $\partial \pi / \partial t = 0$, we have

$$\frac{1}{\beta} \langle \nabla, \mathbf{A}_{\sigma}^{-1} \nabla p(\boldsymbol{x}, t) \rangle + \langle \nabla, p(\boldsymbol{x}, t) \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}) \rangle = 0.$$

which further implies that $\beta^{-1}\nabla\pi + \pi\nabla f(\boldsymbol{x}) = 0$. Solving this equation directly gives $\pi \propto e^{-\beta f(\boldsymbol{x})}$, which completes the proof.

Appendix B. Proof of main theory. In order to bound sampling error between the distribution of the output of LS-SGLD and the target distribution $\pi \propto e^{-\beta f(\boldsymbol{x})}$, we consider a reference sequence generated by LS-LD (2.3), denoted by $\{\boldsymbol{X}_t\}_{t\geq 0}$. Letting $\boldsymbol{X}_0 = \boldsymbol{x}_0$, by the triangle inequality, we can decompose the 2-Wasserstein distance $\mathcal{W}_2(\mathbb{P}(\boldsymbol{x}_K), \pi)$ as follows:

$$\mathcal{W}_2(\mathbb{P}(\boldsymbol{x}_K), \pi) \leq \mathcal{W}_2(\mathbb{P}(\boldsymbol{x}_K), \mathbb{P}(\boldsymbol{X}_{K\eta})) + \mathcal{W}_2(\mathbb{P}(\boldsymbol{X}_{K\eta}), \pi)$$

The first term on the R.H.S. stands for the discretization error of the numerical integrator, and the second term denotes the ergodicity of LS-LD (2.3), which characterizes the mixing time of LS-LD. In what follows, we first deliver the following lemma that characterizes the error term $\mathcal{W}_2(\mathbb{P}(\boldsymbol{X}_{K\eta}), \pi)$.

LEMMA B.1. Under Assumptions 1 and 2, there exists a constant $c_0 \in [||\mathbf{A}_{\sigma}||_2^{-1}, 1]$

$$\mathcal{W}_2(\mathbb{P}(\boldsymbol{X}_{K\eta}), \pi) \le \left[2\lambda \left(\beta f(0) + \log(\Lambda)\right)\right]^{1/2} e^{-c_0 K\eta/(2\beta\lambda)}$$

where $\Lambda = \int_{\mathbb{R}^d} e^{-\beta f(\boldsymbol{x})} d\boldsymbol{x}$ and λ denotes the logarithmic Sobolev constant of the target distribution $\pi \propto e^{-\beta f(\boldsymbol{x})}$.

Note that Lemma B.1 does not require that the target density is log-concave, which can be utilized to prove the convergence rate of LS-SGLD for sampling both log-concave and non-log-concave densities. In the following, we are going to complete the proofs of Theorems 3.3 and 3.7.

B.1. Proof of Theorem 3.3. We first provide the following lemma, which proves an upper bound of the discretization error $W_2(\mathbb{P}(\boldsymbol{x}_K), P(\boldsymbol{X}_{K\eta}))$ for sampling log-concave densities.

LEMMA B.2. Under Assumptions 1, 2, 3, and 4, if we set the step size $\eta \leq Cm\beta^{-1}/M^2$ for some sufficiently small constant C, there exist constants $\gamma_1 \in$

 $[\|\mathbf{A}_{\sigma}\|_{2}^{-2}, 1]$ and $\gamma_{2} = d^{-1} \sum_{i=1}^{d} (1 + 2\sigma - 2\sigma \cos(2\pi i/d))^{-1}$ such that the following holds:

$$\mathcal{W}_2\left(\mathbb{P}(\boldsymbol{x}_K), P(\boldsymbol{X}_{K\eta})\right) \le \left(\frac{2\gamma_1 d\omega^2 K \eta^2}{B}\right)^{1/2} + \left[8\gamma_2 K (K+1)\beta^{-1} d\eta^3\right]^{1/2}$$

Then we can complete the proof of Theorem 3.3 as follows.

Proof of Theorem 3.3. By the triangle inequality and Lemmas B.1 and B.2, it is evident that

$$\mathcal{W}_{2}(\mathbb{P}(\boldsymbol{x}_{K}), \pi) \leq \mathcal{W}_{2}(\mathbb{P}(\boldsymbol{x}_{K}), \mathbb{P}(\boldsymbol{X}_{K\eta})) + \mathcal{W}_{2}(\mathbb{P}(\boldsymbol{X}_{K\eta}), \pi)$$
$$\leq \left(\frac{2\gamma_{1}d\omega^{2}K\eta^{2}}{B}\right)^{1/2} + \left[8\gamma_{2}K(K+1)\beta^{-1}d\eta^{3}\right]^{1/2}$$
$$+ \left[2\lambda\left(\beta f(0) + \log(\Lambda)\right)\right]^{1/2}e^{-c_{0}K\eta/(2\beta\lambda)},$$

which completes the proof.

B.2. Proof of Theorem 3.7. Similar to the proof of Theorem 3.3, we provide the following lemma that characterizes the discretization error $W_2(\mathbb{P}(\boldsymbol{x}_k), \mathbb{P}(\boldsymbol{X}_{k\eta}))$ for sampling from non-log-concave densities.

LEMMA B.3. Under Assumptions 1 and 2, if we set the step size $\eta \leq Cm\beta^{-1}/M^2$ for some sufficiently small constant C, there exist constants $\gamma_1 \in [||\mathbf{A}_{\sigma}||_2^{-2}, 1]$, $\gamma_2 = d^{-1} \sum_{i=1}^d (1 + 2\sigma - 2\sigma \cos(2\pi i/d))^{-1}$, and $\bar{\Gamma} = (3/2 + 2(b + \beta^{-1}d))^{1/2}$ such that the following holds:

$$\mathcal{W}_{2}\left(\mathbb{P}(\boldsymbol{x}_{K}),\mathbb{P}(\boldsymbol{X}_{K\eta})\right) \leq \bar{\Gamma}(K\eta)^{1/2} \left[\left(\frac{\gamma_{1}\beta d\omega^{2}}{2B}K\eta + 2\gamma_{2}\beta M^{2}dK\eta^{2}\right)^{1/2} + \left(\frac{\gamma_{1}\beta d\omega^{2}}{2B}K\eta + 2\gamma_{2}\beta M^{2}dK\eta^{2}\right)^{1/4} \right].$$

Proof of Theorem 3.7. By the triangle inequality and Lemmas B.1 and B.3, it is evident that

$$\begin{aligned} \mathcal{W}_{2}(\mathbb{P}(\boldsymbol{x}_{K}),\pi) &\leq \mathcal{W}_{2}(\mathbb{P}(\boldsymbol{x}_{K}),\mathbb{P}(\boldsymbol{X}_{K\eta})) + \mathcal{W}_{2}(\mathbb{P}(\boldsymbol{X}_{K\eta}),\pi) \\ &\leq \bar{\Gamma}(K\eta)^{1/2} \left[\left(\frac{\gamma_{1}\beta d\omega^{2}}{2B} K\eta + 2\gamma_{2}M^{2} dK\eta^{2} \right)^{1/2} \right. \\ &\left. + \left(\frac{\gamma_{1}\beta d\omega^{2}}{2B} K\eta + 2\gamma_{2}M^{2} dK\eta^{2} \right)^{1/4} \right] \\ &\left. + \left[2\lambda \big(\beta f(0) + \log(\Lambda) \big) \big]^{1/2} e^{-c_{0}K\eta/(2\beta\lambda)}, \end{aligned}$$

which completes the proof.

Appendix C. Proof of lemmas in Appendix B.

C.1. Proof of Lemma B.1. In order to prove Lemma B.1, we require the following lemma.

LEMMA C.1 (Theorem 9.6.1 in [2]). Suppose the target density π satisfies logarithmic Sobolev inequality with a positive constant λ ; for any density μ it holds that

$$\mathcal{W}_2(\mu, \pi) \leq \sqrt{2\lambda D_{KL}(\mu||\pi)}.$$

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

Downloaded 05/07/21 to 155.98.19.70. Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/page/terms

Proof of Lemma B.1. Recall that for LS-LD (2.3), the distribution of X_t , denoted by p(x, t), can be described by the following Fokker–Planck equation:

(C.1)
$$\begin{aligned} \frac{\partial p(\boldsymbol{x},t)}{\partial t} &= \frac{1}{\beta} \langle \mathbf{A}_{\sigma}^{-1}, \nabla^2 p(\boldsymbol{x},t) \rangle + \langle \nabla, p(\boldsymbol{x},t) \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}) \rangle \\ &= \frac{1}{\beta} \langle \nabla, \mathbf{A}_{\sigma}^{-1} \nabla p(\boldsymbol{x},t) \rangle + \langle \nabla, p(\boldsymbol{x},t) \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}) \rangle, \end{aligned}$$

where $\langle \nabla, \mathbf{h}(\boldsymbol{x}) \rangle$ denotes the divergence of the vector field $\mathbf{h}(\boldsymbol{x})$. Let \mathbb{P}_t be the shorthand notation of $p(\boldsymbol{x}, t)$, and denote by $D_{KL}(\mathbb{P}_t || \pi)$ the KL-divergence between the distribution \mathbb{P}_t and the target distribution π . Then, we have

$$\frac{\mathrm{d}D_{KL}(\mathbb{P}_t||\pi)}{\mathrm{d}t} = \int_{\mathbb{R}^d} \frac{\partial}{\partial t} \left[\mathbb{P}_t \log\left(\frac{\mathbb{P}_t}{\pi}\right) \right] \mathrm{d}\boldsymbol{x}$$
$$= \int_{\mathbb{R}^d} \frac{\partial \mathbb{P}_t}{\partial t} \left[\log(\mathbb{P}_t) + 1 - \log(\pi) \right] \mathrm{d}\boldsymbol{x}$$

Similar to the proof of Proposition 2 in [31], by (C.1) we further have

$$\frac{\mathrm{d}D_{KL}(\mathbb{P}_t||\pi)}{\mathrm{d}t} = -\int_{\mathbb{R}^d} \left\langle \frac{1}{\beta} \mathbf{A}_{\sigma}^{-1} \nabla \mathbb{P}_t + \mathbb{P}_t \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}), \nabla \log(\mathbb{P}_t) - \nabla \log(\pi) \right\rangle \mathrm{d}\boldsymbol{x}$$
$$= -\int_{\mathbb{R}^d} \left\langle \mathbf{A}_{\sigma}^{-1} \left(\frac{1}{\beta} \mathbb{P}_t \nabla \log(\mathbb{P}_t) + \mathbb{P}_t \nabla f(\boldsymbol{x}) \right), \nabla \log(\mathbb{P}_t) - \nabla \log(\pi) \right\rangle \mathrm{d}\boldsymbol{x},$$

where the second equality holds due to $\nabla \mathbb{P}_t = \mathbb{P}_t \nabla \log(\mathbb{P}_t)$. In addition, noting that $\pi \propto e^{-\beta f(\boldsymbol{x})}$, we have $\nabla \log(\pi) = -\beta \nabla f(\boldsymbol{x})$. Then we have

$$\frac{\mathrm{d}D_{KL}(\mathbb{P}_t||\pi)}{\mathrm{d}t} = -\frac{1}{\beta} \int_{\mathbb{R}_d} \left\langle \mathbf{A}_{\sigma}^{-1} \big(\nabla \log(\mathbb{P}_t) - \nabla \log(\pi) \big), \nabla \log(\mathbb{P}_t) - \nabla \log(\pi) \big\rangle \mathbb{P}_t \mathrm{d}\mathbf{x} \right\rangle$$
$$= -\frac{1}{\beta} \int_{\mathbb{R}_d} \|\nabla \log(\mathbb{P}_t) - \nabla \log(\pi)\|_{\mathbf{A}_{\sigma}^{-1}}^2 \mathbb{P}_t \mathrm{d}\mathbf{x}.$$

Since \mathbf{A}_{σ} is a positive definite matrix, there exists a constant $c_0 \in [\|\mathbf{A}_{\sigma}\|_2^{-1}, 1]$ such that

(C.2)
$$\frac{\mathrm{d}D_{KL}(\mathbb{P}_t||\pi)}{\mathrm{d}t} \le -\frac{c_0}{\beta} \int_{\mathbb{R}_d} \|\nabla\log(\mathbb{P}_t) - \nabla\log(\pi)\|_2^2 \mathbb{P}_t \mathrm{d}\boldsymbol{x} = -\frac{c_0}{\beta} \mathbf{I}(\mathbb{P}_t||\pi),$$

where $\mathbf{I}(\mathbb{P}_t || \pi)$ denotes the Fisher information between \mathbb{P}_t and π . By Proposition 3.2, we know that the target density π satisfies the logarithmic Sobolev inequality with constant $\lambda > 0$. Then, from [29], we have

$$D_{KL}(\mathbb{P}_t||\pi) \le \frac{1}{\lambda} \mathbf{I}(\mathbb{P}_t||\pi).$$

Plugging the above inequality into (C.2), we obtain

$$\frac{\mathrm{d}D_{KL}(\mathbb{P}_t||\pi)}{\mathrm{d}t} \le -\frac{c_0}{\lambda\beta} D_{KL}(\mathbb{P}_t||\pi),$$

which implies that

$$D_{KL}(\mathbb{P}_t||\pi) \le D_{KL}(\mathbb{P}_0||\pi)e^{-c_0t/(\beta\lambda)}.$$

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

Note that we have $\mathbb{P}_0 = \delta(0)$, where $\delta(\cdot)$ is the Dirac delta function, thus,

$$D_{KL}(\mathbb{P}_0||\pi) = \int_{\mathbb{R}^d} \mathbb{P}_0\left[\log(\mathbb{P}_0) - \log(\pi)\right] \mathrm{d}\boldsymbol{x} = -\log(\pi)|_{\boldsymbol{x}=0} = \beta f(0) + \log(\Lambda),$$

where $\Lambda = \int_{\mathbb{R}^d} e^{-\beta f(\boldsymbol{x})} d\boldsymbol{x}$. Then by Lemma C.1, we have the following regarding the 2-Wasserstein distance $\mathcal{W}_2(\mathbb{P}(\boldsymbol{X}_{k\eta}), \pi)$:

$$\mathcal{W}_2(\mathbb{P}(\boldsymbol{X}_{k\eta}), \pi) \leq \sqrt{2\lambda D_{KL}(\mathbb{P}(\boldsymbol{X}_0)||\pi)} e^{-c_0 t/(2\beta\lambda)}$$
$$= \left[2\lambda (\beta f(0) + \log(\Lambda))\right]^{1/2} e^{-c_0 t/(2\beta\lambda)},$$

which completes the proof.

C.2. Proof of Lemma B.2. We first deliver the following useful lemmas.

LEMMA C.2. Consider any two LS-LD sequences $\{W_t\}_{t\geq 0}$ and $\{V_t\}_{t\geq 0}$, and assume that W_t and V_t have shared Brownian motion terms. Under Assumption 4, for any t > 0 it holds that

$$\mathbb{E}[\|\boldsymbol{W}_t - \boldsymbol{V}_t\|_{\boldsymbol{A}_{\sigma}}^2] \leq \mathbb{E}[\|\boldsymbol{W}_0 - \boldsymbol{V}_0\|_{\boldsymbol{A}_{\sigma}}^2].$$

LEMMA C.3. Under Assumptions 1 and 2, if we set the step size $\eta \leq Cm\beta^{-1}/M^2$ for some sufficiently small constant C, there exists a constant $\gamma_2 = d^{-1} \sum_{i=1}^d (1+2\sigma - 2\sigma \cos(2\pi i/d))^{-1}$ such that for any \boldsymbol{x}_k with $k \geq 0$,

$$\mathbb{E}[\|\mathcal{L}_{\eta}\boldsymbol{x}_{k} - \mathcal{G}_{\eta}\boldsymbol{x}_{k}\|_{\mathbf{A}_{\tau}}^{2}] \leq 4\gamma_{2}\beta^{-1}d\eta^{3}.$$

Now we are ready to complete the proof of Lemma B.2.

Proof of Lemma B.2. For the sake of simplicity, we first define three operators \mathcal{L}_t , \mathcal{G}_t , and \mathcal{S}_t as follows: for any $\boldsymbol{x} \in \mathbb{R}^d$ we denote by $\mathcal{L}_t \boldsymbol{x}$ the random point generated by LS-LD at time t starting from $\boldsymbol{x}, \mathcal{G}_t \boldsymbol{x}$ the point after performing one-step LS-SGLD with full gradient at \boldsymbol{x} with step size t, and $\mathcal{S}_t \boldsymbol{x}$ the point after performing one-step LS-SGLD at \boldsymbol{x} with step size t. Then we have

$$\mathbb{E}[\|\boldsymbol{x}_{K} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}] = \mathbb{E}[\|\boldsymbol{x}_{K} - \mathcal{G}_{\eta}\boldsymbol{x}_{K-1} + \mathcal{G}_{\eta}\boldsymbol{x}_{K-1} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}]$$
(C.3)
$$= \mathbb{E}[\|\boldsymbol{x}_{K} - \mathcal{G}_{\eta}\boldsymbol{x}_{K-1}\|_{\mathbf{A}_{\sigma}}^{2}] + \mathbb{E}[\|\mathcal{G}_{\eta}\boldsymbol{x}_{K-1} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}],$$

where the second equality follows from the fact that $\mathbb{E}[\langle \boldsymbol{x}_{K} - \mathcal{G}_{\eta} \boldsymbol{x}_{K-1}, \mathbf{A}_{\sigma}(\mathcal{G}_{\eta} \boldsymbol{x}_{K-1} - \boldsymbol{X}_{K\eta})\rangle] = 0$ since at any iteration the randomness of the stochastic gradient is independent of the iterate. Regarding the first term on the R.H.S. of (C.3), we have

(C.4)

$$\mathbb{E}[\|\boldsymbol{x}_{K} - \mathcal{G}_{\eta}\boldsymbol{x}_{K-1}\|_{\mathbf{A}_{\sigma}}^{2}] = \eta^{2}\mathbb{E}[\|\boldsymbol{S}_{\eta}\boldsymbol{x}_{K-1} - \mathcal{G}_{\eta}\boldsymbol{x}_{K-1}\|_{\mathbf{A}_{\sigma}}^{2}] \\
\leq \eta^{2}\mathbb{E}[\|\mathbf{A}_{\sigma}^{-1}\mathbf{g}_{K-1} - \mathbf{A}_{\sigma}^{-1}\nabla f(\boldsymbol{x}_{K-1})\|_{\mathbf{A}_{\sigma}}^{2}] \\
\leq \frac{\eta^{2}}{B}\mathbb{E}[\|\mathbf{A}_{\sigma}^{-1}\nabla f_{i}(\boldsymbol{x}_{K-1}) - \mathbf{A}_{\sigma}^{-1}\nabla f(\boldsymbol{x}_{K-1})\|_{\mathbf{A}_{\sigma}}^{2}] \\
\leq \frac{\gamma_{1}\eta^{2}d\omega^{2}}{B},$$

where $\gamma_1 \in [\|\mathbf{A}_{\sigma}\|_2^{-1}, 1)$ is a problem-dependent parameter, the first inequality follows the definitions of operators S_{η} and \mathcal{G}_{η} , the second inequality follows from Lemma A.1 in [25], and the last inequality is by Assumption 3. In terms of the second term on the R.H.S. of (C.3), we have

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

Downloaded 05/07/21 to 155.98.19.70. Redistribution subject to SIAM license or copyright; see https://epubs.siam.org/page/terms

$$\mathbb{E}[\|\mathcal{G}_{\eta}\boldsymbol{x}_{K-1} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}] = \mathbb{E}[\|\mathcal{G}_{\eta}\boldsymbol{x}_{K-1} - \mathcal{L}_{\eta}\boldsymbol{x}_{K-1} + \mathcal{L}_{\eta}\boldsymbol{x}_{K-1} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}]$$

$$\leq (1+\alpha)\mathbb{E}[\|\mathcal{G}_{\eta}\boldsymbol{x}_{K-1} - \mathcal{L}_{\eta}\boldsymbol{x}_{K-1}\|_{\mathbf{A}_{\sigma}}^{2}]$$

$$+ (1+1/\alpha)\mathbb{E}[\|\mathcal{L}_{\eta}\boldsymbol{x}_{K-1} - \mathcal{L}_{\eta}\boldsymbol{X}_{(K-1)\eta}\|_{\mathbf{A}_{\sigma}}^{2}]$$
(C.5)
$$\leq 4(1+\alpha)\gamma_{2}\beta^{-1}d\eta^{3} + (1+1/\alpha)\mathbb{E}[\|\boldsymbol{x}_{K-1} - \boldsymbol{X}_{(K-1)\eta}\|_{\mathbf{A}_{\sigma}}^{2}],$$

where α is a positive constant that will be specified later, the first inequality is by Young's inequality, and the second inequality follows from Lemmas C.2 and C.3. Plugging (C.4) and (C.5) into (C.3) gives

$$\mathbb{E}[\|\boldsymbol{x}_{K} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}] \leq 4(1+\alpha)\gamma_{2}\beta^{-1}d\eta^{3} + \frac{\gamma_{1}\eta^{2}d\omega^{2}}{B} + (1+1/\alpha)\mathbb{E}[\|\boldsymbol{x}_{K-1} - \boldsymbol{X}_{(K-1)\eta}\|_{\mathbf{A}_{\sigma}}^{2}].$$

Then, by recursively applying the above inequality, we obtain

$$\mathbb{E}[\|\boldsymbol{x}_{K} - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^{2}] \leq (1 + 1/\alpha)^{K} \mathbb{E}[\|\boldsymbol{x}_{0} - \boldsymbol{X}_{0}\|_{\mathbf{A}_{\sigma}}^{2}] \\ + \sum_{k=0}^{K-1} (1 + 1/\alpha)^{k} \left[4(1 + \alpha)\gamma_{2}\beta^{-1}d\eta^{3} + \frac{\gamma_{1}\eta^{2}d\omega^{2}}{B} \right] \\ = \alpha \left[(1 + 1/\alpha)^{K} - 1 \right] \cdot \left[4(1 + \alpha)\gamma_{2}\beta^{-1}d\eta^{3} + \frac{\gamma_{1}\eta^{2}d\omega^{2}}{B} \right].$$

Letting $\alpha = K$ and applying the inequality $(1 + 1/K)^K - 1 \le e - 1 \le 2$, the above inequality implies

$$\mathbb{E}[\|\boldsymbol{x}_{K} - \boldsymbol{X}_{K\eta}\|_{\boldsymbol{A}_{\sigma}}^{2}] \leq 2K\eta^{2} \cdot \left[\frac{\gamma_{1}d\omega^{2}}{B} + 4(K+1)\gamma_{2}\beta^{-1}d\eta\right].$$

Based on the definition of 2-Wasserstein distance, we have

$$\mathcal{W}_2^2\big(\mathbb{P}(\boldsymbol{x}_K), \mathbb{P}(\boldsymbol{X}_{K\eta})\big) \le \sqrt{\mathbb{E}[\|\boldsymbol{x}_K - \boldsymbol{X}_{K\eta}\|_2^2]} \le \sqrt{\mathbb{E}[\|\boldsymbol{x}_K - \boldsymbol{X}_{K\eta}\|_{\mathbf{A}_{\sigma}}^2]} \\ \le \left(\frac{2\gamma_1 d\omega^2 K \eta^2}{B}\right)^{1/2} + \left[8\gamma_2 K (K+1)\beta^{-1} d\eta^3\right]^{1/2},$$

where the last inequality is by the fact that $\sqrt{x^2 + y^2} \le |x| + |y|$. This completes the proof.

C.3. Proof of Lemma B.3. In order to prove Lemma B.3, we require the following lemmas.

LEMMA C.4. Under Assumptions 1 and 2, for all $k \ge 0$, there exists a constant $c_1 \in [\|\mathbf{A}_{\sigma}\|_2^{-1}, 1)$ such that

$$\mathbb{E}[\|\boldsymbol{x}_k\|_2^2] \leq \mathbb{E}[\|\boldsymbol{x}_k\|_{\mathbf{A}_{\sigma}}^2] \leq \frac{2(2b+\beta^{-1}d)}{c_1m}.$$

LEMMA C.5 (Theorem 2.3 in [5]). Letting μ, ν be two probability measures with finite exponential second moments, it holds that

$$\mathcal{W}_2(\mu,\nu) \leq \Gamma \left[\sqrt{D_{KL}(\mu||\nu)} + \left[D_{KL}(\mu||\nu) \right]^{1/4} \right],$$

where

$$\Gamma = \inf_{\alpha > 0} \left(\frac{1}{\alpha} \left(\frac{3}{2} + \log \mathbb{E}_{\nu}[e^{\alpha \|\boldsymbol{x}\|_{2}^{2}}] \right) \right)^{1/2}.$$

LEMMA C.6 (Lemma 4 in [38]). Letting $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ be the standard Gaussian random vector with dimension d, it holds that

$$\mathbb{E}\left[\|\mathbf{A}_{\sigma}^{-1}\boldsymbol{\epsilon}\|_{2}^{2}\right] = \sum_{i=1}^{d} \frac{1}{\left(1 + 2\sigma - 2\sigma\cos(2\pi i/d)\right)^{2}}$$

LEMMA C.7. Under Assumptions 1 and 2, let X_t denote the solution of LS-LD (2.3) at time t with initial point $X_0 = 0$. Then if the inverse temperature satisfies $\beta \geq 2 \|\mathbf{A}_{\sigma}\|_2 / m$, it holds that

$$\mathbb{E}[e^{\|\boldsymbol{X}_t\|_2^2}] \le e^{2(b+\beta^{-1}d)t}.$$

Based on the above lemmas, we are able to complete the proof of Lemma B.3.

Proof of Lemma B.3. By Lemma C.5, we know that the 2-Wasserstein distance between any two probability measures can be bounded by their KL divergence. Therefore, the remaining part will focus on deriving the upper bound of the KL divergence $D_{KL}(\mathbb{P}(\boldsymbol{x}_k)||\mathbb{P}(\boldsymbol{X}_{k\eta}))$. Similar to the proof technique used in [15, 35, 40], we leverage the following continuous-time interpolation of LS-SGLD:

(C.6)
$$\tilde{\boldsymbol{X}}_t = \int_0^t -\mathbf{A}_{\sigma}^{-1} \mathbf{G}_s \mathrm{d}s + \int_0^t \sqrt{2\beta^{-1}} \mathbf{A}_{\sigma}^{-1/2} \mathrm{d}\boldsymbol{B}_s,$$

where $\mathbf{G}_t = \sum_{k=0}^{\infty} \mathbf{g}_k \mathbb{1}\{t \in [k\eta, (k+1)\eta)\}$. It can be easily verified that $\tilde{\mathbf{X}}_{k\eta}$ follows the same distribution as \mathbf{x}_k . However, it is worth noting that (C.6) does not form a Markov chain since it contains randomness of the stochastic gradient. To tackle this, we leverage the results in [22] and construct the following Markov chain to mimic (C.6):

$$\hat{\boldsymbol{X}}_t = \int_0^t -\boldsymbol{A}_{\sigma}^{-1} \hat{\boldsymbol{G}}_s ds + \int_0^t \sqrt{2\beta^{-1}} \boldsymbol{A}_{\sigma}^{-1/2} d\boldsymbol{B}_s$$

where $\hat{\mathbf{G}}_s = \mathbb{E}[\mathbf{G}_s | \hat{\mathbf{X}}_s = \hat{\mathbf{X}}_s]$. It was shown that $\hat{\mathbf{X}}_t$ and $\hat{\mathbf{X}}_t$ has the same one-time marginal distribution [22]. Then letting \mathbb{P}_t and \mathbb{Q}_t denote the distribution of \mathbf{X}_t and $\hat{\mathbf{X}}_t$ respectively, by the Girsanov formula, the Radon–Nikodym derivative of \mathbb{P}_t with respect to \mathbb{Q}_t can be derived as follows:

$$\frac{\mathrm{d}\mathbb{P}_t}{\mathrm{d}\mathbb{Q}_t} = \exp\bigg\{\frac{\beta}{2}\int_0^t \langle \nabla f(\hat{\boldsymbol{X}}_s) - \hat{\mathbf{G}}_s, \mathbf{A}_{\sigma}^{-1/2}\mathrm{d}\boldsymbol{B}_s \rangle - \frac{\beta}{4}\int_0^t \|\mathbf{A}_{\sigma}^{-1}\nabla f(\hat{\boldsymbol{X}}_s) - \mathbf{A}_{\sigma}^{-1}\hat{\mathbf{G}}_s\|_2^2\mathrm{d}s\bigg\}.$$

Therefore, letting $T = K\eta$, the KL divergence $D_{KL}(\mathbb{P}_T || \mathbb{Q}_T)$ satisfies

$$D_{KL}(\mathbb{Q}_T || \mathbb{P}_T) = -\int_{\mathbb{R}^d} \log\left(\frac{\mathrm{d}\mathbb{P}_T}{\mathrm{d}\mathbb{Q}_T}\right) \mathrm{d}\mathbb{Q}_T$$

$$= \frac{\beta}{4} \int_0^T \mathbb{E}\left[\|\mathbf{A}_{\sigma}^{-1} \nabla f(\hat{\mathbf{X}}_s) - \mathbf{A}_{\sigma}^{-1} \hat{\mathbf{G}}_s \|_2^2 \right] \mathrm{d}s$$

$$= \frac{\beta}{4} \sum_{k=0}^{K-1} \int_{k\eta}^{(k+1)\eta} \mathbb{E}\left[\|\mathbf{A}_{\sigma}^{-1} \nabla f(\tilde{\mathbf{X}}_s) - \mathbf{A}_{\sigma}^{-1} \mathbf{g}_k \|_2^2 \right] \mathrm{d}s,$$

where the second equality holds due to $\mathbb{E}[\langle \nabla f(\hat{X}_s) - \hat{\mathbf{G}}_s, \mathbf{A}_{\sigma}^{-1/2} \mathrm{d} \boldsymbol{B}_s \rangle] = 0$ and the second equality follows from the fact that \hat{X}_s and \tilde{X}_s follow the same distribution. Using Young's inequality, we have

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

$$\mathbb{E}[\|\mathbf{A}_{\sigma}^{-1}\nabla f(\tilde{\boldsymbol{X}}_{s}) - \mathbf{A}_{\sigma}^{-1}\mathbf{g}_{k}\|_{2}^{2}] \leq 2 \underbrace{\mathbb{E}[\|\mathbf{A}_{\sigma}^{-1}\nabla f(\tilde{\boldsymbol{X}}_{s}) - \mathbf{A}_{\sigma}^{-1}\nabla f(\boldsymbol{x}_{k})\|_{2}^{2}]}_{I_{1}} + 2 \underbrace{\mathbb{E}[\|\mathbf{A}_{\sigma}^{-1}\nabla f(\boldsymbol{x}_{k}) - \mathbf{A}_{\sigma}^{-1}\mathbf{g}_{k}\|_{2}^{2}]}_{I_{2}}.$$

Then we are going to tackle I_1 and I_2 separately. Note that $\|\mathbf{A}_{\sigma}^{-1}\|_2 \leq 1$; thus by Assumption 2, we have the following for I_1 :

$$I_1 \leq \mathbb{E}[\|\nabla f(\tilde{\boldsymbol{X}}_s) - \nabla f(\boldsymbol{x}_k)\|_2^2] \leq M^2 \mathbb{E}[\|\tilde{\boldsymbol{X}}_s - \boldsymbol{x}_k\|_2^2].$$

Based on the definition of $\tilde{\mathbf{X}}_s$, we have $\tilde{\mathbf{X}}_s - \mathbf{x}_k = (s - k\eta) \mathbf{A}_{\sigma}^{-1} \mathbf{g}_k + \sqrt{2\beta^{-1}(s - k\eta)} \mathbf{A}_{\sigma}^{-1/2} \mathbf{\epsilon}_k$. Since $s - k\eta \leq \eta$, it follows that

$$I_{1} \leq M^{2} \mathbb{E}[\|\tilde{\boldsymbol{X}}_{s} - \boldsymbol{x}_{k}\|_{2}^{2}] \leq \eta^{2} M^{2} \mathbb{E}[\|\mathbf{A}_{\sigma}^{-1}\mathbf{g}_{k}\|_{2}^{2}] + 2\eta M^{2} \beta^{-1} \mathbb{E}[\|\mathbf{A}_{\sigma}^{-1/2}\boldsymbol{\epsilon}_{k}\|_{2}^{2}].$$

Regarding I_2 , based on Lemma A.1 in [25] and Assumption 3, we have

$$I_2 \leq \frac{1}{B} \mathbb{E}[\|\mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{x}_k) - \mathbf{A}_{\sigma}^{-1} \nabla f_i(\boldsymbol{x}_k)\|_2^2] \leq \frac{\gamma_1 d\omega^2}{B},$$

where $\gamma_1 \in [\|\mathbf{A}_{\sigma}\|_2^{-2}, 1)$ is a problem-dependent parameter. Putting everything together, we have

$$D_{KL}(\mathbb{Q}_T||\mathbb{P}_T) \leq \sum_{k=0}^{K-1} \eta \bigg\{ \frac{\beta}{2} \eta^2 M^2 \mathbb{E} \big[\|\mathbf{g}_k\|_{\mathbf{A}_{\sigma}^{-2}}^2 \big] + \eta M^2 \mathbb{E} \big[\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}_{\sigma}^{-1}}^2 \big] + \frac{\gamma_1 \beta d\omega^2}{2B} \bigg\}.$$

By Lemma D.1 and Young's inequality, we know that

$$\mathbb{E}\left[\|\mathbf{g}_k\|_{\mathbf{A}_{\sigma}^{-2}}^2\right] \le \mathbb{E}[\|\mathbf{g}_k\|_2^2] \le 2M^2 \mathbb{E}[\|\mathbf{x}_k\|_2^2] + 2G^2 \le \frac{4M^2(2b + \beta^{-1}d)}{c_1 m} + 2G^2,$$

where $G = \max_{i \in [n]} \|\nabla f_i(0)\|_2$. Then by Lemma C.6, we have $\mathbb{E}[\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}_{\sigma}^{-1}}^2] \leq \gamma_2 d$ with $\gamma_2 = d^{-1} \sum_{i=1}^d (1 + 2\sigma - 2\sigma \cos(2\pi i/d))^{-1}$, which is strictly smaller than 1. Therefore,

$$D_{KL}(\mathbb{Q}_T || \mathbb{P}_T) \le \frac{2\beta M^4(2b + \beta^{-1}d) + \beta c_1 m M^2 G^2}{c_1 m} K\eta^3 + \gamma_2 M^2 dK \eta^2 + \frac{\gamma_1 \beta d\omega^2}{2B} K\eta.$$

For sufficiently small step size such that

$$\eta \le \frac{c_1 \beta^{-1} \gamma_2 m d}{2M^2 (2b + \beta^{-1} d) + c_1 m G^2},$$

we have

$$D_{KL}(\mathbb{P}(\boldsymbol{x}_K)||\mathbb{P}(\boldsymbol{X}_{K\eta})) \leq \frac{\gamma_1 \beta d\omega^2}{2B} K\eta + 2\gamma_2 M^2 dK\eta^2.$$

Then, by Lemma C.5, we have

$$\mathcal{W}_2\big(\mathbb{P}(\boldsymbol{x}_K),\mathbb{P}(\boldsymbol{X}_{K\eta})\big) \leq \Gamma\Big[\sqrt{D_{KL}\big(\mathbb{P}(\boldsymbol{x}_K)||\mathbb{P}(\boldsymbol{X}_{K\eta})\big)} + \big[D_{KL}\big(\mathbb{P}(\boldsymbol{x}_K)||\mathbb{P}(\boldsymbol{X}_{K\eta})\big)\big]^{1/4}\Big],$$

where Γ can be further bounded as

$$\Gamma \leq \left(\frac{3}{2} + \log \mathbb{E}[e^{\|\boldsymbol{X}_{K_{\eta}}\|_{2}^{2}}]\right)^{1/2}$$

$$\leq \left(3/2 + 2(b + \beta d)K\eta\right)^{1/2}$$

$$\leq \left(3/2 + 2(b + \beta d)\right)^{1/2} \cdot (K\eta)^{1/2},$$

where the first inequality is by the choice $\alpha = 1$, the second inequality is by Lemma C.7, and the last inequality is by the assumption that $K\eta > 1$. Therefore, defining $\overline{\Gamma} = (3/2 + 2(b + \beta^{-1}d))^{1/2}$, the 2-Wasserstein distance $\mathcal{W}_2(\mathbb{P}(\boldsymbol{x}_K), \mathbb{P}(\boldsymbol{X}_{K\eta}))$ can be bounded by

$$\mathcal{W}_{2}\left(\mathbb{P}(\boldsymbol{x}_{K}),\mathbb{P}(\boldsymbol{X}_{K\eta})\right) \leq \bar{\Gamma}(K\eta)^{1/2} \left[\left(\frac{\gamma_{1}\beta d\omega^{2}}{2B}K\eta + 2\gamma_{2}M^{2}dK\eta^{2}\right)^{1/2} + \left(\frac{\gamma_{1}\beta d\omega^{2}}{2B}K\eta + 2\gamma_{2}M^{2}dK\eta^{2}\right)^{1/4} \right],$$

which completes the proof.

Appendix D. Proof of lemmas in Appendix C.

D.1. Proof of Lemma C.2.

Proof of Lemma C.2. Assuming shared Brownian motions in W_t and V_t , we have

$$d\mathbb{E}[\|\boldsymbol{W}_{t} - \boldsymbol{V}_{t}\|_{\boldsymbol{A}_{\sigma}}^{2}] = -2\mathbb{E}[\langle \boldsymbol{A}_{\sigma}^{-1} (\nabla f(\boldsymbol{W}_{t}) - \nabla f(\boldsymbol{V}_{t})), \boldsymbol{A}_{\sigma}(\boldsymbol{W}_{t} - \boldsymbol{V}_{t}) \rangle] dt$$

$$= -2\mathbb{E}[\langle \nabla f(\boldsymbol{W}_{t}) - \nabla f(\boldsymbol{V}_{t}), \boldsymbol{W}_{t} - \boldsymbol{V}_{t} \rangle] dt$$

$$\leq 0,$$

where the first equality follows from the fact that we assume shared Brownian motion terms on both dynamics $\{W_t\}_{t\geq 0}$ and $\{V_t\}_{t\geq 0}$ and the inequality is due to the convexity of f(x). Therefore, it can be evidently concluded that

$$\mathbb{E}[\|\boldsymbol{W}_t - \boldsymbol{V}_t\|_{\boldsymbol{A}_{\sigma}}^2] \le \mathbb{E}[\|\boldsymbol{W}_0 - \boldsymbol{V}_0\|_{\boldsymbol{A}_{\sigma}}^2],$$
oof.

which completes the proof.

D.2. Proof of Lemma C.3.

Proof of Lemma C.3. To simplify the analysis, let \boldsymbol{x} be any iterate of LS-SGLD and define $\boldsymbol{x} = \boldsymbol{X}_0$. Then the operators \mathcal{G}_{η} and \mathcal{L}_{η} satisfy

$$\mathcal{G}_{\eta}\boldsymbol{x} = \boldsymbol{X}_{0} - \eta \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{X}_{0}) + \sqrt{2\beta^{-1}\eta} \mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}$$

$$= \boldsymbol{X}_{0} - \int_{0}^{\eta} \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{X}_{0}) dt + \int_{0}^{\eta} \sqrt{2\beta^{-1}} \mathbf{A}_{\sigma}^{-1/2} d\boldsymbol{B}_{t};$$

$$\mathcal{L}_{\eta}\boldsymbol{x} = \boldsymbol{X}_{0} - \int_{0}^{\eta} \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{X}_{t}) dt + \int_{0}^{\eta} \sqrt{2\beta^{-1}} \mathbf{A}_{\sigma}^{-1/2} d\boldsymbol{B}_{t}.$$

Considering synchronous Brownian terms in \mathcal{G}_{η} and \mathcal{L}_{η} , we have

$$\mathbb{E}[\|\mathcal{L}_{\eta}\boldsymbol{x} - \mathcal{G}_{\eta}\boldsymbol{x}\|_{\mathbf{A}_{\sigma}}^{2}] = \mathbb{E}\left[\left\|\int_{0}^{\eta} \left[\mathbf{A}_{\sigma}^{-1}\nabla f(\boldsymbol{X}_{0}) - \mathbf{A}_{\sigma}^{-1}\nabla f(\boldsymbol{X}_{t})\right] \mathrm{d}t\right\|_{\mathbf{A}_{\sigma}}^{2}\right]$$
$$\leq \mathbb{E}\left[\eta\int_{0}^{\eta} \left\|\mathbf{A}_{\sigma}^{-1}\left[\nabla f(\boldsymbol{X}_{0}) - \nabla f(\boldsymbol{X}_{t})\right]\right\|_{\mathbf{A}_{\sigma}}^{2} \mathrm{d}t\right]$$
$$\leq M^{2}\left[\eta\int_{0}^{\eta} \mathbb{E}[\|\boldsymbol{X}_{t} - \boldsymbol{X}_{0}\|_{2}^{2}] \mathrm{d}t\right],$$

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

where the second inequality follows from Jensen's inequality and the last inequality follows from Assumption 2 and the fact that $\|\mathbf{A}_{\sigma}\|_{2} \geq 1$. We further have

$$\mathbb{E}[\|\boldsymbol{X}_t - \boldsymbol{X}_0\|_2^2] = \mathbb{E}\left[\left\|\int_0^t \mathbf{A}_{\sigma}^{-1} \nabla f(\boldsymbol{X}_{\tau}) \mathrm{d}\tau\right\|_2^2\right] + 2t\beta^{-1} \mathbb{E}[\|\mathbf{A}_{\sigma}^{-1/2}\boldsymbol{\epsilon}\|_2^2]$$
$$\leq \mathbb{E}\left[t\int_0^t \|\nabla f(\boldsymbol{X}_{\tau})\|_2^2 \mathrm{d}\tau\right] + 2\gamma_2 t\beta^{-1} d,$$

where the inequality is by Jensen's inequality and Lemma C.6 and $\gamma_2 = d^{-1} \sum_{i=1}^{d} (1 + 2\sigma - 2\sigma \cos(2\pi i/d))^{-1}$ is strictly smaller than 1. By Lemma D.1, we have

$$\mathbb{E}[\|\nabla f(\boldsymbol{X}_{\tau})\|_{2}^{2}] \leq 2M^{2}\mathbb{E}[\|\boldsymbol{X}_{\tau}\|_{2}^{2}] + 2G^{2}$$

Note that by Ito's lemma we have for any $0 \le s \le \tau$,

$$\frac{\mathrm{d}\mathbb{E}[\|\boldsymbol{X}_s\|_{\boldsymbol{A}_{\sigma}}^2]}{\mathrm{d}s} = -2\mathbb{E}[\langle \boldsymbol{X}_s, \nabla f(\boldsymbol{X}_s)\rangle] + \beta^{-1}d \le -2m\mathbb{E}[\|\boldsymbol{X}_s\|_2^2] + 2b + \beta^{-1}d \le 2b + \beta^{-1}d,$$

where the second inequality follows from Assumption 1. Therefore,

$$\mathbb{E}[\|\boldsymbol{X}_{\tau}\|_{2}^{2}] \leq \mathbb{E}[\|\boldsymbol{X}_{\tau}\|_{\boldsymbol{A}_{\sigma}}^{2}] = \mathbb{E}[\|\boldsymbol{X}_{0}\|_{\boldsymbol{A}_{\sigma}}^{2}] + \int_{0}^{\tau} \frac{\mathrm{d}\mathbb{E}[\|\boldsymbol{X}_{s}\|_{\boldsymbol{A}_{\sigma}}^{2}]}{\mathrm{d}s} \mathrm{d}s \leq \mathbb{E}[\|\boldsymbol{X}_{0}\|_{\boldsymbol{A}_{\sigma}}^{2}] + \tau(2b + \beta^{-1}d).$$

Note that $X_0 = x$ is an iterate of LS-SGLD; by Lemma C.4 we have $\mathbb{E}[||X_0||^2_{\mathbf{A}_{\sigma}}] \leq (2b + \beta^{-1}d)/(c_1m)$ for some constant $c_1 \in [||\mathbf{A}_{\sigma}||^{-1}_2, 1]$. Therefore,

$$\mathbb{E}[\|\nabla f(\boldsymbol{X}_{\tau})\|_{2}^{2}] \leq 2M^{2}\mathbb{E}[\|\boldsymbol{X}_{\tau}\|_{2}^{2}] + 2G^{2} \leq \frac{4M^{2}(2b+\beta^{-1}d)}{c_{1}m} + 2G^{2} + 2M^{2}\tau(2b+\beta^{-1}d).$$

Thus, it follows that

$$\mathbb{E}[\|\boldsymbol{X}_t - \boldsymbol{X}_0\|_2^2] \le \left(\frac{4M^2(2b + \beta^{-1}d)}{c_1m} + 2G^2 + 2M^2\tau(2b + \beta^{-1}d)\right)t^2 + 2\gamma_2 t\beta^{-1}d.$$

Noting that $\tau, t \leq \eta$, plugging the above inequality into (D.1), we have

$$\mathbb{E}[\|\mathcal{L}_{\eta}\boldsymbol{x} - \mathcal{G}_{\eta}\boldsymbol{x}\|_{2}^{2}] \leq M^{2} \bigg[\bigg(\frac{4M^{2}(2b+\beta^{-1}d)}{c_{1}m} + 2G^{2} + 2M^{2}(2b+\beta^{-1}d)\eta \bigg) \eta^{4} + 2\gamma_{2}\beta^{-1}d\eta^{3} \bigg].$$

For sufficiently small step size satisfying

$$\eta \leq \frac{c_1 \beta^{-1} \gamma_2 m d}{4M^2 (2b + \beta^{-1} d) + 2c_1 m G^2} \wedge \sqrt{\frac{\gamma_2 \beta^{-1} d}{M^2 (2b + \beta^{-1} d)}},$$

where \wedge stands for maximum of two numbers,

we have

$$\mathbb{E}[\|\mathcal{L}_{\eta}\boldsymbol{x} - \mathcal{G}_{\eta}\boldsymbol{x}\|_{2}^{2}] \leq 4\gamma_{2}\beta^{-1}d\eta^{3}.$$

This completes the proof.

D.3. Proof of Lemma C.4. In order to prove Lemma C.4, we need the following lemma.

LEMMA D.1 (Lemma 3.1 in [35]). For any $\boldsymbol{x} \in \mathbb{R}^d$ and $i \in [n]$, it holds that

$$\|\nabla f_i(\boldsymbol{x})\|_2 \le M \|\boldsymbol{x}\|_2 + G,$$

where $G = \max_{i \in [n]} \|\nabla f_i(0)\|_2$.

Proof of Lemma C.4. Recall the update formula of x_k ,

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \eta \mathbf{A}_{\sigma}^{-1} \mathbf{g}_k + \sqrt{2\beta^{-1}\eta} \mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}_k.$$

Therefore, it holds that

A50

$$\mathbb{E}[\|\boldsymbol{x}_{k+1}\|_{\mathbf{A}_{\sigma}}^{2}] = \mathbb{E}[\|\boldsymbol{x}_{k} - \eta \mathbf{A}_{\sigma}^{-1} \mathbf{g}_{k}\|_{\mathbf{A}_{\sigma}}^{2}] + 2\eta\beta^{-1}\mathbb{E}[\|\mathbf{A}_{\sigma}^{-1/2} \boldsymbol{\epsilon}_{k}\|_{\mathbf{A}_{\sigma}}^{2}] \\ = \mathbb{E}[\|\boldsymbol{x}_{k}\|_{\mathbf{A}_{\sigma}}^{2}] - 2\eta\mathbb{E}[\langle \boldsymbol{x}_{k}, \mathbf{g}_{k} \rangle] + \eta^{2}\mathbb{E}[\|\mathbf{g}_{k}\|_{\mathbf{A}_{\sigma}^{-1}}^{2}] + 2\eta\beta^{-1}d,$$

where the second equality follows from the fact that $\mathbb{E}[\|\boldsymbol{\epsilon}_k\|_2^2] = d$. Noting that all eigenvalues of \mathbf{A}_{σ} are greater than 1, it follows that

$$\begin{split} \mathbb{E}[\|\boldsymbol{x}_{k+1}\|_{\mathbf{A}_{\sigma}}^{2}] &= \mathbb{E}[\|\boldsymbol{x}_{k}\|_{\mathbf{A}_{\sigma}}^{2}] - 2\eta \mathbb{E}[\langle \boldsymbol{x}_{k}, \nabla f(\boldsymbol{x}) \rangle] + \eta^{2} \mathbb{E}[\|\mathbf{g}_{k}\|_{2}^{2}] + 2\eta\beta^{-1}d \\ &\leq \mathbb{E}[\|\boldsymbol{x}_{k}\|_{\mathbf{A}_{\sigma}}^{2}] - 2\eta m \mathbb{E}[\|\boldsymbol{x}_{k}\|_{2}^{2}] + 2\eta b + 2\eta^{2}(M^{2} \mathbb{E}[\|\boldsymbol{x}_{k}\|_{2}^{2}] + G^{2}) + 2\eta\beta^{-1}d, \end{split}$$

where the inequality follows from Assumption 1, Lemma D.1, and Young's inequality. Since the step size η satisfies $\eta \leq m/(2M^2)$, we further have

$$\mathbb{E}[\|\boldsymbol{x}_{k+1}\|_{\mathbf{A}_{\sigma}}^2] \leq \mathbb{E}[\|\boldsymbol{x}_k\|_{\mathbf{A}_{\sigma}}^2] - \eta m \mathbb{E}[\|\boldsymbol{x}_k\|_2^2] + 2\eta (b + \beta^{-1}d + \eta G^2).$$

Recall that all eigenvalues of \mathbf{A}_{σ} are greater than 1; there exists a constant $\|\mathbf{A}_{\sigma}\|_{2}^{-1} \leq c_{1} \leq 1$ such that

(D.2)
$$\mathbb{E}[\|\boldsymbol{x}_{k+1}\|_{\mathbf{A}_{\sigma}}^{2}] \leq (1 - c_{1}\eta m)\mathbb{E}[\|\boldsymbol{x}_{k}\|_{\mathbf{A}_{\sigma}}^{2}] + 2\eta(b + \beta^{-1}d + \eta G^{2}).$$

Since $\eta \leq 1/(c_1m) \wedge b/G$, (D.2) implies that the following holds for all $k \geq 0$:

$$\mathbb{E}[\|\boldsymbol{x}_k\|_{\mathbf{A}_{\sigma}}^2] \le (1 - c_1 \eta m)^k \|\boldsymbol{x}_0\|_{\mathbf{A}_{\sigma}}^2 + \frac{2(2b + \beta^{-1}d)}{c_1 m}.$$

Since at the initialization $\boldsymbol{x}_0 = 0$, we have

$$\mathbb{E}[\|\boldsymbol{x}_k\|_2^2] \leq \mathbb{E}[\|\boldsymbol{x}_k\|_{\mathbf{A}_{\sigma}}^2] \leq \frac{2(2b+\beta^{-1}d)}{c_1m}.$$

This completes the proof.

D.4. Proof of Lemma C.7.

Proof of Lemma C.7. We first define the function $L(t) = e^{\|\mathbf{X}_t\|_{\mathbf{A}_{\sigma}}^2}$; then by Ito's formula, we have

$$d\mathbb{E}[L(t)] = -2\mathbb{E}[\langle \mathbf{A}_{\sigma} \mathbf{X}_{t}, \mathbf{A}_{\sigma}^{-1} \nabla f(\mathbf{X}_{t}) \rangle L(t)] dt + \mathbb{E}[\langle 4\mathbf{A}_{\sigma} \mathbf{X}_{t} \mathbf{X}_{t}^{\top} \mathbf{A}_{\sigma} + 2\mathbf{A}_{\sigma}, \beta^{-1} \mathbf{A}_{\sigma}^{-1} \mathbf{I} \rangle L(t)] dt = -2\mathbb{E}[\langle (\mathbf{X}_{t}, \nabla f(\mathbf{X}_{t}) \rangle - \beta^{-1} d - 2\beta^{-1} \| \mathbf{X}_{t} \|_{\mathbf{A}_{\sigma}}^{2}) L(t)] dt.$$

By Assumption 1, we further have

$$d\mathbb{E}[L(t)] \le 2\mathbb{E}\left[\left((-m\|\boldsymbol{X}_t\|_2^2 + 2\beta^{-1}\|\boldsymbol{X}_t\|_{\boldsymbol{A}_{\sigma}}^2) + b + \beta^{-1}d\right)L(t)\right]dt.$$

Therefore, assuming $\beta \geq 2 \|\mathbf{A}_{\sigma}\|_2 / m$, we have

$$d\mathbb{E}[L(t)] \le 2(b + \beta^{-1}d)\mathbb{E}[L(t)]dt.$$

Since L(t) is always positive, it holds that

$$\mathbb{E}[L(t)] \le L(0)e^{2(b+\beta^{-1}d)t}.$$

Noting that $\|\boldsymbol{X}_t\|_{\boldsymbol{A}_{\sigma}}^2 \geq \|\boldsymbol{X}_t\|_2^2$, we immediately have

$$\mathbb{E}[e^{\|\boldsymbol{X}_t\|_2^2}] \le \mathbb{E}[e^{\|\boldsymbol{X}_t\|_{\mathbf{A}_{\sigma}}^2}] \le L(0)e^{2(b+\beta^{-1}d)t},$$

which completes the proof.

REFERENCES

- A. ATTIA, R. ŞTEFĂNESCU, AND A. SANDU, The reduced-order hybrid Monte Carlo sampling smoother, Internat. J. Numer. Methods Fluids, 83 (2017), pp. 28-51.
- [2] D. BAKRY, I. GENTIL, AND M. LEDOUX, Analysis and Geometry of Markov Diffusion Operators, Grundlehren Math. Wiss. 348, Springer, Berlin, 2013.
- [3] R. BARDENET, A. DOUCET, AND C. HOLMES, On Markov chain Monte Carlo methods for tall data, J. Mach. Learn. Res., 18 (2017), pp. 1515–1557.
- M. BETANCOURT, The fundamental incompatibility of scalable hamiltonian Monte Carlo and naive data subsampling, in Proceedings of the International Conference on Machine Learning, 2015, pp. 533–540.
- [5] F. BOLLEY AND C. VILLANI, Weighted Csiszár-Kullback-Pinsker inequalities and applications to transportation inequalities, Ann. Fac. Sci. Toulouse Math. (6), 14 (2005), https://doi. org/10.5802/afst.1095.
- [6] J. BORGGAARD, N. GLATT-HOLTZ, AND J. KROMETIS, A Bayesian approach to estimating background flows from a passive scalar, SIAM/ASA J. Uncertain. Quantif., 8 (2020), pp. 1036– 1060.
- [7] A. BOVIER, M. ECKHOFF, V. GAYRARD, AND M. KLEIN, Metastability in reversible diffusion processes I: Sharp asymptotics for capacities and exit times, J. Eur. Math. Soc. (JEMS), 6 (2004), pp. 399–424.
- [8] 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, preprint, arXiv:1802.05431, 2018.
- C. CHEN, N. DING, AND L. CARIN, On the convergence of stochastic gradient MCMC algorithms with high-order integrators, in Advances in Neural Information Processing Systems, 2015, pp. 2278–2286.
- [10] T. CHEN, E. FOX, AND C. GUESTRIN, Stochastic gradient Hamiltonian Monte Carlo, in Proceedings of the International Conference on Machine Learning, 2014.
- [11] Y. CHEN, R. DWIVEDI, M. J. WAINWRIGHT, AND B. YU, Fast Mixing of Metropolized Hamiltonian Monte Carlo: Benefits of Multi-Step Gradients, preprint, arXiv:1905.12247, 2019.
- [12] X. CHENG, N. S. CHATTERJI, Y. ABBASI-YADKORI, P. L. BARTLETT, AND M. I. JORDAN, Sharp Convergence Rates for Langevin Dynamics in the Nonconvex Setting, preprint, arXiv:1805.01648, 2018.
- [13] X. CHENG, N. S. CHATTERJI, P. L. BARTLETT, AND M. I. JORDAN, Underdamped Langevin mcmc: A non-asymptotic analysis, in Proceedings of the 31st Conference on Learning Theory, Vol. 75, 2018, pp. 300–323.
- [14] T.-S. CHIANG, C.-R. HWANG, AND S. J. SHEU, Diffusion for global optimization in rⁿ, SIAM J. Control Optim., 25 (1987), pp. 737–753.
- [15] 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.

- [16] A. S. DALALYAN AND A. G. KARAGULYAN, User-Friendly Guarantees for the Langevin Monte Carlo with Inaccurate Gradient, preprint, arXiv:1710.00095, 2017.
- [17] K.-D. DANG, M. QUIROZ, R. KOHN, M.-N. TRAN, AND M. VILLANI, Hamiltonian Monte Carlo with energy conserving subsampling, J. Mach. Learn. Res., 20 (2019), pp. 1–31.
- [18] K. A. DUBEY, S. J. REDDI, S. A. WILLIAMSON, B. POCZOS, 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.
- [19] A. DURMUS AND E. MOULINES, Nonasymptotic convergence analysis for the unadjusted Langevin algorithm, Ann. Appl. Probab., 27 (2017), pp. 1551–1587.
- [20] R. DWIVEDI, Y. CHEN, M. J. WAINWRIGHT, AND B. YU, Log-concave sampling: Metropolis-Hastings algorithms are fast!, in Proceedings of the 31st Conference on Learning Theory, 2018, pp. 793–797.
- [21] G. GOLUB AND C. VAN LOAN, Matrix Computation, 3rd ed., Johns Hopkins University Press, Baltimore, 1996.
- [22] I. GYÖNGY, Mimicking the one-dimensional marginal distributions of processes having an Itô differential, Probab. Theory Related fields, 71 (1986), pp. 501–516.
- [23] P. E. KLOEDEN AND E. PLATEN, Numerical Solution of Stochastic Differential Equations, Springer, Berlin, 1992.
- [24] L. KREUSSER, S. OSHER, AND B. WANG, A Deterministic Approach to Avoid Saddle Points, preprint, arXiv:1901.06827, 2019.
- [25] L. LEI, C. JU, J. CHEN, AND M. I. JORDAN, Non-convex finite-sum optimization via SCSG methods, in Advances in Neural Information Processing Systems 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds., 2017, pp. 2348–2358.
- [26] C. LI, C. CHEN, D. CARLSON, AND C. LAWRENCE, Preconditioned stochastic gradient Langevin dynamics for deep neural networks, in Proceedings of the Association for the Advancement of Artificial Intelligence, 2016.
- [27] Z. LI, T. ZHANG, AND J. LI, Stochastic Gradient Hamiltonian Monte Carlo with Variance Reduction for Bayesian Inference, preprint, arXiv:1803.11159, 2018.
- [28] 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.
- [29] P. A. MARKOWICH AND C. VILLANI, On the trend to equilibrium for the Fokker-Planck equation: An interplay between physics and functional analysis, in Physics and Functional Analysis, Matematica Contemporanea 19, Brazilian Mathematical Society, 1999.
- [30] J. C. MATTINGLY, A. M. STUART, AND D. J. HIGHAM, Ergodicity for SDEs and approximations: Locally Lipschitz vector fields and degenerate noise, Stochastic Process. Appl., 101 (2002), pp. 185–232.
- [31] W. MOU, L. WANG, X. ZHAI, AND K. ZHENG, Generalization Bounds of SGLD for Non-Convex Learning: Two Theoretical Viewpoints, preprint, arXiv:1707.05947, 2017.
- [32] R. M. NEAL, MCMC using hamiltonian dynamics, in Handbook of Markov Chain Monte Carlo, Handbooks Modern Statist. Methods 2, Routledge, 2011, pp. 113–162.
- [33] S. OSHER, B. WANG, P. YIN, X. LUO, M. PHAM, AND A. LIN, Laplacian Smoothing Gradient Descent, preprint, arXiv:1806.06317, 2018.
- [34] G. PARISI, Correlation functions and computer simulations, Nuclear Phys. B, 180 (1981), pp. 378–384.
- [35] M. RAGINSKY, A. RAKHLIN, AND M. TELGARSKY, Non-convex learning via stochastic gradient Langevin dynamics: A nonasymptotic analysis, in Proceedings of the Conference on Learning Theory, 2017, pp. 1674–1703.
- [36] 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.
- [37] S. J. VOLLMER, K. C. ZYGALAKIS, AND Y. W. TEH, Exploration of the (non-) asymptotic bias and variance of stochastic gradient Langevin dynamics, J. Mach. Learn. Res., 17 (2016), pp. 5504–5548.
- [38] B. WANG, Q. GU, M. BOEDIHARDJO, F. BAREKAT, AND S. OSHER, DP-LSSGD: A Stochastic Optimization Method to Lift the Utility in Privacy-Preserving ERM, preprint, arXiv:1906.12056, 2019.
- [39] M. WELLING AND T. Y. WHYE, *Bayesian learning via stochastic gradient Langevin dynamics*, in Proceedings of the International Conference on Machine Learning, 2011.
- [40] P. XU, J. CHEN, D. ZOU, AND Q. GU, Global convergence of Langevin dynamics based algorithms for nonconvex optimization, in Advances in Neural Information Processing Systems, 2018, pp. 3122–3133.

A53

- [41] C. ZHANG, B. SHAHBABA, AND H. ZHAO, Hamiltonian Monte Carlo acceleration using surrogate functions with random bases, Statist. Comput., 27 (2017), pp. 1473–1490.
- [42] D. ZOU, P. XU, AND Q. GU, Stochastic variance-reduced Hamilton Monte Carlo methods, in Proceedings of the 35th International Conference on Machine Learning, 2018, pp. 6028– 6037.
- [43] 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.
- [44] D. ZOU, P. XU, AND Q. GU, Sampling from non-log-concave distributions via variance-reduced gradient Langevin dynamics, in Artificial Intelligence and Statistics, Proc. Machine Learning Res. 89, MLR Press, 2019, pp. 2936–2945.
- [45] D. ZOU, P. XU, AND Q. GU, Stochastic gradient Hamiltonian Monte Carlo methods with recursive variance reduction, in Advances in Neural Information Processing Systems, 2019.