
Improved Convergence Rate of Stochastic Gradient Langevin Dynamics with Variance Reduction and its Application to Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The stochastic gradient Langevin Dynamics is one of the most fundamental al-
2 gorithms to solve sampling problems and non-convex optimization appearing in
3 several machine learning applications. Especially, its variance reduced versions
4 have nowadays gained particular attention. In this paper, we study two variants
5 of this kind, namely, the Stochastic Variance Reduced Gradient Langevin Dynam-
6 ics and the Stochastic Recursive Gradient Langevin Dynamics. We prove their
7 convergence to the objective distribution in terms of KL-divergence under the
8 sole assumptions of smoothness and Log-Sobolev inequality which are weaker
9 conditions than those used in prior works for these algorithms. With the batch
10 size and the inner loop length set to \sqrt{n} , the gradient complexity to achieve an
11 ϵ -precision is $\tilde{O}((n + dn^{1/2}\epsilon^{-1})\gamma^2L^2\alpha^{-2})$, which is an improvement from any
12 previous analyses. We also show some essential applications of our result to
13 non-convex optimization.

14 1 Introduction

15 1.1 Background and Organization

16 Over the past decade, the gradient Langevin Dynamics (GLD) has gained particular attention for
17 providing an effective tool for sampling from a Gibbs distribution, a fundamental task omnipresent
18 in the field of machine learning and statistics, and for non-convex optimization, which is nowadays
19 witnessing an unignorable empirical success. Notably, GLD is a stochastic differential equation (SDE)
20 that can be viewed as the steepest descent flow of the Kullback-Leibler (KL) divergence towards the
21 stationary Gibbs distribution in the space of measures endowed with the 2-Wasserstein metric (Jordan
22 et al., 1998). As a consequence of the unique properties of GLD, its implementable discrete schemes
23 and their ability to suitably track it have been the subject of a large number of studies.

24 The Euler-Maruyama scheme of GLD gives rise to an algorithm known as the Langevin Monte
25 Carlo method (LMC). This algorithm is biased (Wibisono, 2018): that is, the distribution of the
26 discrete scheme does not converge to the same as GLD. Nonetheless, it has been shown that this
27 bias could be made arbitrarily small under certain assumptions by taking a sufficiently small step
28 size (Dalalyan, 2017b; Vempala and Wibisono, 2019). Dalalyan (2017a,b) provided one of the first
29 non-asymptotic rates of convergence of LMC for smooth log-concave distributions. Assumptions to
30 obtain a non-asymptotic analysis and this controllable bias have been relaxed by further research to
31 dissipativity and smoothness (Raginsky et al., 2017; Xu et al., 2018), and recently to Log-Sobolev
32 inequality (LSI) and smoothness (Vempala and Wibisono, 2019). This relaxation of conditions
33 is especially meaningful as the objective distribution nowadays tends to become more and more
34 complicated beyond the classical assumption of log-concavity.

35 However, in the field of machine learning, the main function can often be formulated as the average
36 of the loss function of an enormous number of training data points (Welling and Teh, 2011), which
37 subsequently makes it difficult to calculate its full gradient. As a result, research on stochastic
38 algorithms has been also conducted to avoid this computational burden (Chen et al., 2021; Dubey
39 et al., 2016; Raginsky et al., 2017; Welling and Teh, 2011; Xu et al., 2018; Zou et al., 2018, 2019a,b,
40 2021). Welling and Teh (2011) introduced the concept of Stochastic Gradient Langevin Dynamics
41 (SGLD) which combines the Stochastic Gradient Descent with LMC. This has been the subject of
42 successful studies (Raginsky et al., 2017; Welling and Teh, 2011; Xu et al., 2018). Nevertheless, the
43 variance of its stochastic gradient is too large, which leads to a slow convergence compared to LMC.
44 Therefore, stochastic gradient Langevin Dynamics algorithms with variance reduction, such as the
45 Stochastic Variance Reduced Gradient Langevin Dynamics (SVRG-LD), have been considered and
46 their convergence has been thoroughly analyzed for both sampling (Dubey et al., 2016; Zou et al.,
47 2018, 2019a, 2021) and optimization (Huang and Becker, 2021; Xu et al., 2018).

48 Dubey et al. (2016) first united SGLD with variance reduction techniques and proposed two new
49 algorithms, namely, SVRG-LD and SAGA-LD. Chatterji et al. (2018) and Zou et al. (2018) proved
50 the convergence rate of SVRG-LD to the target distribution in 2-Wasserstein distance for smooth
51 log-concave distributions. Xu et al. (2018) showed the weak convergence of SVRG-LD under the
52 smoothness and dissipativity conditions. They expanded the non-asymptotic analysis of Raginsky
53 et al. (2017) to LMC and SVRG-LD, and improved the result for SGLD. Few years ago, Zou et al.
54 (2019a) provided the gradient complexity of SVRG-LD to converge to the stationary distribution in
55 2-Wasserstein distance under the smoothness and dissipativity assumptions. This convergence can be
56 even improved if we make a warm-start (Zou et al., 2021). While these works investigated algorithms
57 with fixed hyperparameters, Huang and Becker (2021) additionally assumed a strict saddle and some
58 other minor conditions to study SVRG-LD with a decreasing step size and improved its convergence
59 in high probability to the second order stationary point. Zou et al. (2019b) also applied variance
60 reduction techniques to the Hamiltonian Langevin Dynamics, or underdamped Langevin Dynamics
61 in opposition to GLD also known as overdamped Langevin Dynamics. As we can observe, the
62 current convergence analyses of the stochastic schemes with variance reduction are mostly restricted
63 to log-concavity and dissipativity, and do not enjoy the same broad convergence guarantee with a
64 concrete gradient complexity as LMC does under LSI and smoothness in terms of KL-divergence.

65 Therefore, in order to bridge this theoretical gap between LMC and stochastic gradient Langevin
66 Dynamics with variance reduction, we study in this paper the convergence of the latter under the
67 relaxed assumptions of smoothness and LSI. In Section 3, we study the convergence to the Gibbs
68 distribution of SVRG-LD and the Stochastic Recursive Gradient Langevin Dynamics (SARAH-LD),
69 another variant of stochastic gradient Langevin Dynamics with variance reduction inspired by the
70 Stochastic Recursive Gradient algorithm (SARAH) of Nguyen et al. (2017a,b). On the other hand,
71 optimization and sampling are only two sides of the same coin for GLD. That is why, in Section 4,
72 we also investigate implications of Section 3 for non-convex optimization. We prove the convergence
73 of SVRG-LD and SARAH-LD to the global minimum of dissipative functions and we provide their
74 non-asymptotic rate of convergence. We also consider the additional weak Morse assumption and
75 study its effect.

76 1.2 Contributions

77 The major contributions of this paper can be summarized as follows. We provide a non-asymptotic
78 analysis of the convergence of SVRG-LD and SARAH-LD to the Gibbs distribution in terms of
79 KL-divergence under smoothness and LSI which are weaker conditions than those used in prior works
80 for these algorithms. KL-divergence is generally a stronger convergence criterion than both total
81 variation (TV) and 2-Wasserstein distance as they can be controlled by KL-divergence under the LSI
82 condition. Notably, we prove that, with the batch size and inner loop length set to \sqrt{n} , the gradient
83 complexity to achieve an ϵ -precision in terms of KL-divergence is $\tilde{O}((n + dn^{1/2}\epsilon^{-1})\gamma^2 L^2 \alpha^{-2})$,
84 which is better than any previous analyses. See Table 1 for a comparison with previous research in
85 terms of assumptions, criterion and gradient complexity. We also prove the convergence of SVRG-
86 LD and SARAH-LD to the global minimum under an additional assumption of dissipativity with a
87 gradient complexity of $\tilde{O}((n + n^{1/2}\epsilon^{-1}dL\alpha^{-1})\gamma^2 L^2 \alpha^{-2})$ which is better than previous work since
88 it has almost all the time a dependence on n of $O(\sqrt{n})$ and does not require the batch size and the
89 inner loop length to depend on the accuracy ϵ . On the other hand, we import the idea of Li and
90 Erdogdu (2020) from product manifolds of spheres to the Euclidean space in order to show that under

Table 1: Comparison of our main result with prior works (sampling). The first three works are about LMC. Compared to Vempala et al. (2019), with the same assumptions and criterion, the order of gradient complexity is improved from n to \sqrt{n} . The others are about SVRG-LD except the last one which is about the Stochastic Gradient Hamiltonian Monte Carlo Methods with Recursive Variance Reduction. ϵ is the accuracy required on the criterion, d is the dimension of the input of the main function, n is the number of data points, and L is the smoothness constant. * 2-Wass. stands for “2-Wasserstein”, and conv. stands for “convergence”. ** $\text{poly}(M, L)$ stands for a polynomial of M and L .

Method	Major Assumptions	Criterion*	Gradient Complexity**
Dalalyan (2017a)	Smooth, Log-concave (M)	2-Wass.	$\tilde{O}\left(\frac{nd}{\epsilon^2} \cdot \text{poly}(M, L)\right)$
Xu et al. (2018)	Smooth, Dissipative	Weak conv.	$\tilde{O}\left(\frac{nd}{\epsilon}\right) \cdot e^{\tilde{O}(d)}$
Vempala et al. (2019)	Smooth, Log-Sobolev (α)	KL	$\tilde{O}\left(\frac{n}{\epsilon} \cdot d\gamma^2 L^2 \alpha^{-2}\right)$
Zou et al. (2018)	Smooth, Log-concave (M)	2-Wass.	$\tilde{O}\left(n + \frac{L^{3/2} n^{1/2} d^{1/2}}{M^{3/2} \epsilon}\right)$
Zou et al. (2019a)	Smooth, Dissipative	2-Wass.	$\tilde{O}\left(n + \frac{n^{3/4}}{\epsilon^2} + \frac{n^{1/2}}{\epsilon^4}\right) \cdot e^{\tilde{O}(\gamma+d)}$
Zou et al. (2021)	Smooth, Dissipative, Warm-start	TV	$\tilde{O}\left(\frac{\gamma^2}{\epsilon^2}\right) \cdot e^{\tilde{O}(d)}$
Zou et al. (2019b)	Smooth, Dissipative	2-Wass.	$\tilde{O}\left(\left(n + \frac{n^{1/2}}{\epsilon^2 \mu_*^{3/2}}\right) \wedge \frac{\mu_*^{-2}}{\epsilon^4}\right)$
This paper	Smooth, Log-Sobolev (α)	KL	$\tilde{O}\left(\left(n + \frac{dn^{1/2}}{\epsilon}\right) \cdot \gamma^2 L^2 \alpha^{-2}\right)$

91 the additional assumption of weak Morse, the convergence in the Euclidean space can be accelerated
 92 by eliminating the exponential dependence on $1/\epsilon$.

93 1.3 Other Related Works

94 The theoretical study of GLD goes back to Chiang et al. (1987) who showed that global convergence
 95 could be achieved with a proper annealing schedule. This work did not specify how to implement
 96 this SDE, but Gelfand and Mitter (1991) filled this gap. Later, Borkar and Mitter (1999) proved an
 97 asymptotic convergence in terms of relative entropy for the discrete scheme of gradient Langevin
 98 Dynamics when the inverse temperature and the step size are kept constant.

99 The variance reduction technique, introduced to Langevin Dynamics by Dubey et al. (2016), was origi-
 100 nally presented by Johnson and Zhang (2013) as Stochastic Variance Reduced Gradient (SVRG) to
 101 improve the convergence speed of Stochastic Gradient Descent. Other variance reduction techniques
 102 were also considered such as the Stochastic Recursive Gradient Langevin Dynamics (SARAH) from
 103 Nguyen et al. (2017a,b) which outperforms SVRG in non-convex optimization (Pham et al., 2020)
 104 and is used in many algorithms such as SSRGD (Li, 2019) and SpiderBoost (Wang et al., 2019).

105 Li and Erdogdu (2020) extended Vempala and Wibisono’s result to Riemannian manifolds. One of
 106 the highlights of their work is that they showed the Log-Sobolev constant of the Gibbs distribution
 107 for a product manifold of spheres only depends on a polynomial of the inverse temperature under
 108 some particular conditions including weak Morse. We will adapt this result to our situation.

109 In the concurrent work of Balasubramanian et al. (2022) (especially Section 6), they also studied
 110 the convergence of stochastic schemes of GLD with more relaxed conditions than prior analyses.
 111 However, our contributions are not overshadowed by theirs, and we clarify the reasons. In Subsection
 112 6.1 of their paper, Balasubramanian et al. (2022) focused on stochastic discrete schemes with finite
 113 variance and bias (which is not the case for SVRG-LD) and provided a first-order convergence
 114 guarantee in the space of measures equipped with the 2-Wasserstein distance. Subsection 6.2 proved
 115 a global convergence under some other conditions but most of these two analyses did not consider
 116 in particular the usual case in machine learning when F is the average of some other functions,
 117 which leads to a generally worse gradient complexity than ours. Concerning this finite sum setting,
 118 Balasubramanian et al. (2022) investigated the Variance Reduced LMC algorithm (slightly different
 119 from SVRG-LD in this paper) in Subsection 6.3 and gave a first-order convergence under the sole
 120 assumption of smoothness. When restrained in our problem setting, the gradient complexity of
 121 SVRG-LD and SARAH-LD we provide is still considerably better (see Section 3 for more details).

122 **1.4 Notation**

123 We denote deterministic vectors by a lower case symbol (e.g., x) and random variables by an upper
 124 case symbol (e.g., X). The Euclidean norm is denoted by $\|\cdot\|$ for vectors and the inner product
 125 by $\langle \cdot, \cdot \rangle$. For matrices, $\|\cdot\|$ is the norm induced by the Euclidean norm for vectors. We only
 126 treat distributions absolutely continuous with respect to the Lebesgue measure in \mathbb{R}^d for simplicity.
 127 Especially, throughout the paper, ν refers to the probability measure with the density function
 128 $d\nu \propto e^{-\gamma F} dx$, where F is a function introduced below. $a \vee b$ is equivalent to $\max\{a, b\}$ and $a \wedge b$
 129 to $\min\{a, b\}$. We also use the shorthand \tilde{O} to hide logarithmic polynomials.

130 **2 Preliminaries**

131 In this section, we briefly explain the problem setting, necessary mathematical background and
 132 assumptions used in this paper.

133 **2.1 Problem Setting and GLD**

134 In Section 3, we consider sampling from a distribution written in the form $d\nu \propto e^{-\gamma F} dx$ where γ is
 135 a positive constant (which corresponds to the inverse temperature) and $F: \mathbb{R}^d \rightarrow \mathbb{R}$ is formulated as
 136 $F(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$, the average of the loss function of n training data points $\{x^{(i)}\}_{i=1}^n$. Here,
 137 $f_i(x) := f(x, x^{(i)})$ can be regarded as the loss of data $x^{(i)}$. For instance, F can be the average
 138 of the negative log likelihood of n training data points. In Section 4, we consider the non-convex
 139 optimization (minimization) of the same F as above.

140 GLD can be described as the following stochastic differential equation (SDE):

$$dX_t^{\text{GLD}} = -\nabla F(X_t^{\text{GLD}})dt + \sqrt{2/\gamma}dB(t), \quad (1)$$

141 where $\gamma > 0$ is called the inverse temperature parameter and $\{B(t)\}_{t \geq 0}$ is the standard Brownian
 142 motion in \mathbb{R}^d . It can be used for sampling since under some reasonable assumptions of F , the
 143 distribution ρ_t^{GLD} of X_t^{GLD} governed by SDE (1) converges to the invariant stationary distribution
 144 $d\nu \propto e^{-\gamma F} dx$, also known as the Gibbs distribution (Chiang et al., 1987). Moreover, as previously
 145 mentioned, this convergence is efficient in the sense that SDE (1) corresponds to the steepest descent
 146 flow of the Kullback-Leibler (KL) divergence towards the stationary distribution in the space of
 147 measures endowed with the 2-Wasserstein metric (Jordan et al., 1998). Alternatively, GLD can be
 148 interpreted as the composite optimization problem of a negative entropy and an expected function
 149 value as follows (Wibisono, 2018):

$$\min_{q: \text{density}} \mathbb{E}_q[\gamma F] + \mathbb{E}_q[\log q].$$

150 The gradient flow is the well-known Fokker-Planck equation associated to SDE (1):

$$\frac{\partial \rho_t^{\text{GLD}}}{\partial t} = \nabla \cdot (\rho_t^{\text{GLD}} \nabla F) + \frac{1}{\gamma} \Delta \rho_t^{\text{GLD}} = \frac{1}{\gamma} \nabla \cdot \left(\rho_t^{\text{GLD}} \nabla \log \frac{\rho_t^{\text{GLD}}}{\nu} \right). \quad (2)$$

151 This will be useful in our analysis. In addition to its potential for sampling, GLD can also be employed
 152 for non-convex optimization as the Gibbs distribution concentrates on the global minimum of F for
 153 sufficiently large values of γ (Hwang, 1980).

154 **2.2 Algorithms of GLD**

155 Applying the Euler-Maruyama scheme to (1), we obtain the Langevin Monte Carlo (LMC)

$$X_{k+1} = X_k - \eta \nabla F(X_k) + \sqrt{2\eta/\gamma} \epsilon_k,$$

156 where η is called the step size. This is similar to the gradient descent except the additional Gaussian
 157 noise $\sqrt{2\eta/\gamma} \epsilon_k$, where $\epsilon_k \sim N(0, I_{d \times d})$ and $I_{d \times d}$ is the $d \times d$ unit matrix. In the case n is huge and
 158 the computation of ∇F is too difficult, we are incited to use stochastic gradient methods in analogy
 159 to stochastic gradient optimization. This gives

$$X_{k+1} = X_k - \eta v(X_k) + \sqrt{2\eta/\gamma} \epsilon_k,$$

Algorithm 1: SVRG-LD / SARAH-LD

```
1 input: step size  $\eta > 0$ , batch size  $B$ , epoch length  $m$ , inverse temperature  $\gamma \geq 1$ 
2 initialization:  $X_0 = 0$ ,  $X^{(0)} = X_0$ 
3 foreach  $s = 0, 1, \dots, (K/m)$  do
4    $v_{sm} = \nabla F(X^{(s)})$ 
5   randomly draw  $\epsilon_{sm} \sim N(0, I_{d \times d})$ 
6    $X_{sm+1} = X_{sm} - \eta v_{sm} + \sqrt{2\eta/\gamma} \epsilon_{sm}$ 
7   foreach  $l = 1, \dots, m - 1$  do
8      $k = sm + l$ 
9     randomly pick a subset  $I_k$  from  $\{1, \dots, n\}$  of size  $|I_k| = B$ 
10    randomly draw  $\epsilon_k \sim N(0, I_{d \times d})$ 
11    if SVRG-LD then
12       $v_k = \frac{1}{B} \sum_{i_k \in I_k} (\nabla f_{i_k}(X_k) - \nabla f_{i_k}(X^{(s)})) + v_{sm}$ 
13    else if SARAH-LD then
14       $v_k = \frac{1}{B} \sum_{i_k \in I_k} (\nabla f_{i_k}(X_k) - \nabla f_{i_k}(X_{k-1})) + v_{k-1}$ 
15    end
16     $X_{k+1} = X_k - \eta v_k + \sqrt{2\eta/\gamma} \epsilon_k$ 
17  end
18   $X^{(s+1)} = X_{(s+1)m}$ 
19 end
```

160 where $v(X_k)$ is the stochastic gradient. When $v(X_k)$ is defined as $\frac{1}{B} \sum_{i_k \in I_k} \nabla f_{i_k}(X_k)$, where
161 B is called the batch size and I_k is a random subset uniformly chosen from $\{1, \dots, n\}$ such
162 that $|I_k| = B$, we obtain the Stochastic Gradient Langevin Dynamics (SGLD). As this method
163 exhibits a slow convergence, it has been popular to use variance reduction methods such as
164 the Stochastic Variance Reduced Gradient Langevin Dynamics (SVRG-LD) where $v(X_k) =$
165 $\frac{1}{B} \sum_{i_k \in I_k} (\nabla f_{i_k}(X_k) - \nabla f_{i_k}(X^{(s)})) + \nabla F(X^{(s)})$. Details of this algorithm is stated in Algo-
166 rithm 1. $X^{(s)}$ is a reference point updated every m steps so that $X_{sm} = X^{(s)}$. As we can observe in
167 Lemma A.4, around the optimal point, the variance of the stochastic gradient is indeed decreased
168 as $X^{(s)}$ and X_k are both close to each other. We can also easily extend some successful stochastic
169 gradient algorithms to Langevin Dynamics. Hence, we are motivated to extend the Stochastic Recur-
170 sive Gradient Algorithm (SARAH) to Langevin Dynamics since we can expect that some bottlenecks
171 of the analysis of SVRG-LD can be removed in that of SARAH-LD as subtracting the previous
172 stochastic gradient enables a stabler performance than SVRG-LD. This algorithm can be described as
173 Algorithm 1 with $v(X_k) = \frac{1}{B} \sum_{i_k \in I_k} (\nabla f_{i_k}(X_k) - \nabla f_{i_k}(X_{k-1})) + v(X_{k-1})$.

174 **Definition 1.** We define ρ_k as the distribution of X_k generated at the k th step of SVRG-LD, and
175 similarly ϕ_k for SARAH-LD.

176 2.3 Assumptions

177 The assumptions used throughout this paper can be summarized as follows.

178 **Assumption 1.** For all $i = 1, \dots, n$, ∇f_i is twice differentiable, and $\forall x, y \in \mathbb{R}^d$, $\|\nabla^2 f_i(x)\| \leq L$.
179 In other words, f_i ($i = 1, \dots, n$) and F are L -smooth.

180 **Assumption 2.** Distribution ν satisfies the Log-Sobolev inequality (LSI) with a constant α . That is,
181 for all probability density functions ρ absolutely continuous with respect to ν , the following holds:

$$H_\nu(\rho) \leq \frac{1}{2\alpha} J_\nu(\rho),$$

182 where $H_\nu(\rho) := \mathbb{E}_\rho \left[\log \frac{\rho}{\nu} \right]$ is the KL-divergence of ρ with respect to ν , and $J_\nu(\rho) :=$
183 $\mathbb{E}_\rho \left[\|\nabla \log \frac{\rho}{\nu}\|^2 \right]$ is the relative Fisher information of ρ with respect to ν .

184 The recent work of Vempala and Wibisono (2019) motivates us to use the combination of smoothness
185 and LSI for the analysis of SVRG-LD and SARAH-LD. Indeed, they showed that these conditions

186 were enough to assure for the Euler-Maruyama scheme an exponentially fast convergence and a bias
187 controllable by the step size. Under smoothness, LSI is not only the necessary condition of log-
188 concavity and dissipativity, but is also robust to bounded perturbation and Lipschitz mapping, contrary
189 to log-concavity (Vempala and Wibisono, 2019). For example, for any distribution $d\nu$ that satisfies
190 LSI and bounded function $B : \mathbb{R}^d \rightarrow \mathbb{R}$, $d\tilde{\nu} \propto e^B d\nu$ satisfies LSI as well (Holley and Stroock, 1986).
191 Moreover, while KL-divergence is not in general convex with regard to the Wasserstein geodesic,
192 thanks to LSI, the Polyak-Łojaciewicz condition is satisfied. It is well-known that LSI suffices to
193 realize an exponential convergence for the case of continuous time Langevin Dynamics (Vempala
194 and Wibisono, 2019). That is why, it is actually both useful and natural to suppose LSI in this context.
195 Note that under L -smoothness of F and LSI with constant α for $d\nu \propto e^{-\gamma F} dx$, it holds that $\alpha \leq \gamma L$
196 (Vempala and Wibisono, 2019).

197 As for optimization, we additionally use the following conditions.

198 **Assumption 3.** F is (M, b) -dissipative. That is, there exist constants $M > 0$ and $b > 0$ such that for
199 all $x \in \mathbb{R}^d$ the following holds: $\langle \nabla F(x), x \rangle \geq M\|x\|^2 - b$.

200 **Assumption 4** (Li and Erdogdu (2020), Assumption 3.3). F satisfies the weak Morse condition. That
201 is, for all non-zero eigenvalues of the Hessian of stationary points, there exists a constant $\lambda^\dagger \in (0, 1]$
202 such that

$$\lambda^\dagger \leq \inf \{ |\lambda_i(\nabla^2 F(x))| \mid \nabla F(x) = 0, i \in 1, \dots, d, \lambda_i(\nabla^2 F(x)) \neq 0 \}.$$

203 Furthermore, for the set \mathcal{S} of stationary points that are not a global minimum,
204 $\sup_{x \in \mathcal{S}} \lambda_{\min}(\nabla^2 F(x)) \leq -\lambda^\dagger$.

205 **Assumption 5.** $\nabla^2 f_i$ is L' -Lipschitz and without loss of generality, we let $\min_{x \in \mathbb{R}^d} F(x) = 0$.

206 **Assumption 6.** F has a unique global minimum.

207 Smoothness and dissipativity are a classical combination of assumptions for this kind of problem
208 setting (Raginsky et al., 2017; Xu et al., 2018; Zou et al., 2019a). We assume dissipativity instead of
209 LSI for non-convex optimization in order to obtain an explicit value of the Log-Sobolev constant
210 of $d\nu \propto e^{-\gamma F} dx$ in function of the inverse temperature parameter γ (see Property C.3), making a
211 non-asymptotic analysis possible. Furthermore, Assumptions 4 to 6 can ameliorate the exponential
212 dependence of the inverse of the Log-Sobolev constant on the inverse temperature parameter to a
213 polynomial one (see Property C.4).

214 3 Main Results

215 In this section, we state our main results which prove that SVRG-LD and SARAH-LD (Algorithm 1)
216 achieve an exponentially fast convergence to the Gibbs distribution and a controllable bias in terms
217 of KL-divergence under the sole assumptions of LSI and smoothness. We provide their gradient
218 complexity as well. The proofs can be found in Appendix A and B respectively.

219 3.1 Improved Convergence of SVRG-LD

220 Our analysis shows that the convergence of SVRG-LD to the stationary distribution $d\nu \propto e^{-\gamma F} dx$
221 can be formulated as the theorem below.

222 **Theorem 1.** Under Assumptions 1 and 2, $0 < \eta < \frac{\alpha}{16\sqrt{6}L^2 m \gamma}$, $\gamma \geq 1$ and $B \geq m$, for all
223 $k = 1, 2, \dots$, the following holds in the update of SVRG-LD where $\Xi = \frac{(n-B)}{B(n-1)}$:

$$H_\nu(\rho_k) \leq e^{-\frac{\alpha \eta}{\gamma} k} H_\nu(\rho_0) + \frac{224\eta\gamma dL^2}{3\alpha} (2 + 3\Xi + 2m\Xi).$$

224 We observe that the bias term of the upper bound, which is the second term linearly dependent on η ,
225 can be easily controlled while the first term exponentially converges to 0 with $k \rightarrow \infty$. This is more
226 precisely formulated in the following corollary.

227 **Corollary 1.1.** Under the same assumptions as Theorem 1, for all $\epsilon \geq 0$, if we choose step size
228 η such that $\eta \leq \frac{3\alpha\epsilon}{448\gamma dL^2}$, then a precision $H_\nu(\rho_k) \leq \epsilon$ is reached after $k \geq \frac{\gamma}{\alpha\eta} \log \frac{2H_\nu(\rho_0)}{\epsilon}$ steps.

229 Especially, if we take $B = m = \sqrt{n}$ and the largest permissible step size $\eta = \frac{\alpha}{16\sqrt{6}L^2\sqrt{n}\gamma} \wedge \frac{3\alpha\epsilon}{448dL^2\gamma}$,
 230 then the gradient complexity becomes

$$\tilde{O}\left(\left(n + \frac{dn^{\frac{1}{2}}}{\epsilon}\right) \cdot \frac{\gamma^2 L^2}{\alpha^2}\right).$$

231 This gradient complexity is an improvement compared with prior works for three reasons. First of all,
 232 we provide a non-asymptotic analysis of the convergence of SVRG-LD under smoothness and Log-
 233 Sobolev inequality which are conditions weaker than those (e.g., log-concavity or dissipativity) used in
 234 prior works for these algorithms. Moreover, we prove it in terms of KL-divergence which is generally
 235 a stronger convergence criterion than both total variation (TV) and 2-Wasserstein distance as they can
 236 both be controlled by KL-divergence under the LSI condition. For instance, TV was used by [Zou et al.](#)
 237 [\(2021\)](#) and 2-Wasserstein distance by [Dalalyan \(2017a\)](#) and [Zou et al. \(2019a\)](#). KL-divergence makes
 238 it possible to unify these two different criteria. Finally, while prior research generally used Girsanov's
 239 theorem which generates a bias term that accumulates through the iteration (see for example [Raginsky](#)
 240 [et al. \(2017\)](#) and [Xu et al. \(2018\)](#)), we solve this issue by taking benefit of the exponential convergence
 241 of GLD to the Gibbs distribution under LSI and smoothness that enables us to forget about past
 242 bias. That way, with the batch size and inner loop set to \sqrt{n} , the gradient complexity to achieve
 243 an ϵ -precision in terms of KL-divergence becomes $\tilde{O}((n + dn^{1/2}\epsilon^{-1})\gamma^2 L^2 \alpha^{-2})$, which is better
 244 than previous analyses. For example, [Vempala and Wibisono \(2019\)](#) provided a gradient complexity
 245 of $\tilde{O}(n\epsilon^{-1}d\gamma^2 L^2 \alpha^{-2})$ for LMC under Assumptions 1 and 2, and [Zou et al. \(2019a\)](#) a gradient
 246 complexity of $\tilde{O}(n + n^{3/4}\epsilon^{-2} + n^{1/2}\epsilon^{-4}) \cdot e^{\tilde{O}(\gamma+d)}$ for SVRG-LD under Assumptions 1 and 3. Note
 247 that the dependence on the dimension d is not improved since α^{-1} may exponentially depend on d .
 248 Recently, [Zou et al. \(2019b\)](#) proposed the Stochastic Gradient Hamiltonian Monte Carlo Methods
 249 with Recursive Variance Reduction with a gradient complexity of $\tilde{O}((n + n^{1/2}\epsilon^{-2}\mu_*^{-3/2}) \wedge \mu_*^{-2}\epsilon^{-4})$
 250 in terms of 2-Wasserstein distance. Even though their algorithm is based on the underdamped
 251 Langevin Dynamics whose discrete schemes use to perform better than those of the overdamped
 252 Langevin Dynamics such as SVRG-LD, our gradient complexity, which applies to a broader family
 253 of distributions, is almost the same except for a small interval of ϵ , but we do not require the batch
 254 size B and the inner loop length m to depend on ϵ while [Zou et al. \(2019b\)](#) do, i.e., $B \lesssim B_0^{1/2}$,
 255 $m = O(B_0/B)$, where $B_0 = \tilde{O}(\epsilon^{-4}\mu_*^{-1} \wedge n)$. This strengthens the importance of our result since
 256 it shows that adapting this analysis to other stochastic schemes of GLD is promising and could lead
 257 to tighter bounds and relaxation of conditions. See Table 1 for a summary. Concerning the concurrent
 258 work of [Balasubramanian et al. \(2022\)](#), under the sole assumption of smoothness, they provided
 259 a gradient complexity of $O(L^2 d^2 n / \epsilon^2)$ for the Variance Reduced LMC algorithm that updates the
 260 stochastic gradient differently as SVRG-LD and SARAH-LD. This is almost the square of our result,
 261 and in some extent, our work can be interpreted as an acceleration of their result with a slightly
 262 stronger additional condition than Poincaré inequality.

263 **Proof Sketch** Proceeding in a similar way as [Vempala and Wibisono \(2019\)](#), we evaluate how
 264 $H_\nu(\rho_k)$ decreases at each step as shown in Theorem A.1 of Appendix A. This is realized by comparing
 265 the evolution of the continuous-time GLD for time η and one step of SVRG-LD. Since we use a
 266 stochastic gradient, we need at the same time to evaluate the variance of the stochastic gradient.
 267 Theorem 1 can be obtained by recursively solving the inequality derived in Theorem A.1.

268 3.2 Convergence Analysis of SARAH-LD

269 As for SARAH-LD, its convergence to the stationary distribution $d\nu \propto e^{-\gamma F} dx$ can be formulated
 270 as the theorem below. Interestingly, we obtain the same result as SVRG-LD (Theorem 1) but we do
 271 not require $B \geq m$ anymore.

272 **Theorem 2.** Under Assumptions 1 and 2, $0 < \eta < \frac{\alpha}{16\sqrt{2}L^2 m \gamma}$ and $\gamma \geq 1$, for all $k = 1, 2, \dots$, the
 273 following holds in the update of SARAH-LD where $\Xi = \frac{(n-B)}{B(n-1)}$:

$$H_\nu(\phi_k) \leq e^{-\frac{\alpha\eta}{\gamma}k} H_\nu(\phi_0) + \frac{32\eta\gamma dL^2}{3\alpha} (2 + \Xi + 2m\Xi).$$

274 This is the first convergence guarantee of SARAH-LD in this problem setting so far, and it leads to
 275 the following gradient complexity.

276 **Corollary 2.1.** *Under the same assumptions as Theorem 2, for all $\epsilon \geq 0$, if we choose step*
 277 *size η such that $\eta \leq \frac{3\alpha\epsilon}{64\gamma dL^2} (2 + \Xi + 2m\Xi)^{-1}$, then a precision $H_\nu(\phi_k) \leq \epsilon$ is reached after*
 278 *$k \geq \frac{\gamma}{\alpha\eta} \log \frac{2H_\nu(\phi_0)}{\epsilon}$ steps. Especially, if we take $B = m = \sqrt{n}$ and the largest permissible step size*
 279 *$\eta = \frac{\alpha}{16\sqrt{2}L^2\sqrt{n}\gamma} \wedge \frac{3\alpha\epsilon}{320dL^2\gamma}$, then the gradient complexity becomes*

$$\tilde{O} \left(\left(n + \frac{dn^{\frac{1}{2}}}{\epsilon} \right) \cdot \frac{\gamma^2 L^2}{\alpha^2} \right).$$

280 The reason why we obtain the same gradient complexity for both SARAH-LD and SVRG-LD (except
 281 better coefficients for SARAH-LD) is that in our analysis, the Brownian noise added at each step
 282 of the Langevin Dynamics plays the role of a fundamental bottleneck that even SARAH-LD could
 283 not eliminate, and we still need to set $B = m = \sqrt{n}$. We can hypothesize that this order of gradient
 284 complexity might be tight for variance-reduced stochastic gradient Langevin Dynamics algorithms.

285 4 Some Applications to Non-Convex Optimization

286 Here, we apply our main results to non-convex optimization. Thanks to our analysis applicable to a
 287 broader family of probability distributions satisfying LSI, the additional conditions we pose in this
 288 section are mainly reflected in the concrete formulation of the Log-Sobolev constant, which keeps our
 289 study simple and clear. The proofs can be found in Appendix C. Since SVRG-LD and SARAH-LD
 290 exhibited almost the same performance in sampling, we can simultaneously analyse them. We first
 291 prove the convergence to the global minimum of SVRG-LD and SARAH-LD without clarifying the
 292 explicit formulation of the Log-Sobolev constant in function of γ .

293 **Theorem 3.** *Using SVRG-LD or SARAH-LD, under Assumptions 1 to 3, $0 < \eta < \frac{\alpha}{16\sqrt{6}L^2m\gamma}$,*
 294 *$\gamma \geq \frac{4d}{\epsilon} \log \left(\frac{eL}{M} \right) \vee \frac{8db}{\epsilon^2} \vee 1 \vee \frac{2}{M}$ and $B \geq m$, if we take $B = m = \sqrt{n}$ and the largest permissible*
 295 *step size $\eta = \frac{\alpha}{16\sqrt{6}L^2\sqrt{n}\gamma} \wedge \frac{3}{1792} \frac{\alpha^2\epsilon}{L^2d\gamma}$, the gradient complexity to reach a precision of*

$$\mathbb{E}_{X_k}[F(X_k)] - F(X^*) \leq \epsilon$$

296 *is*

$$\tilde{O} \left(\left(n + \frac{n^{\frac{1}{2}}}{\epsilon} \cdot \frac{dL}{\alpha} \right) \frac{\gamma^2 L^2}{\alpha^2} \right),$$

297 *where α is a function of γ , and X^* is the global minimum of F .*

298 **Remark 1.** *Under Assumptions 1 and 3, Assumption 2 is negligible as shown in Property C.2.*

299 Under Assumptions 1 to 3 only, this leads to a gradient complexity which exponentially depends on
 300 the inverse of the precision level ϵ as shown in the next corollary since the inverse of the Log-Sobolev
 301 constant exponentially depends on γ .

302 **Corollary 3.1.** *Under the same assumptions as Theorem 3, taking $\gamma = i(\epsilon) := \frac{4d}{\epsilon} \log \left(\frac{eL}{M} \right) \vee \frac{8db}{\epsilon^2} \vee$
 303 $1 \vee \frac{2}{M}$, we obtain a gradient complexity of*

$$\tilde{O} \left(\left(n + \frac{n^{\frac{1}{2}}}{\epsilon} \cdot \frac{dL}{C_1 i(\epsilon)} e^{C_2 i(\epsilon)} \right) L^2 e^{2C_2 i(\epsilon)} \right)$$

304 *since $\alpha = \gamma C_1 e^{-C_2 \gamma}$ (Property C.3).*

305 The second term with $n^{1/2}$ is almost all the time dominant since it has a factor that exponentially
 306 depends on $1/\epsilon$ and the first term not. This dependence on n of $O(n^{1/2})$ is the best so far for these
 307 algorithms. Moreover, comparing with the gradient complexity $\tilde{O}(n^{1/2} \lambda^{-4} \epsilon^{-5/2}) \cdot e^{\tilde{O}(d)}$, also of
 308 order $n^{1/2}$, provided by Xu et al. (2018) who used SVRG-LD and the same assumptions, our gradient
 309 complexity is an improvement since their analysis required a batch size B and an inner loop length
 310 m that strongly depend on ϵ (i.e., $B = \sqrt{n}\epsilon^{-3/2}$, $m = \sqrt{n}\epsilon^{3/2}$) and ours does not. Note that the
 311 dependence of the gradient complexity of Xu et al. (2018) on $1/\epsilon$ is not necessarily better than ours
 312 as λ is actually the spectral gap of the discrete-time Markov chain generated by (1) and its inverse

313 exponentially depends on $1/\epsilon$ as well. Although Xu et al. (2018) did not investigate the explicit
 314 nature of λ , this is supported by Raginsky et al. (2017) who proved this exponential dependence for
 315 the spectral gap of the continuous-time SDE and by Mattingly et al. (2002) who showed the spectral
 316 gap of continuous-time SDE and that of discrete-time version are almost the same in this context.

317 **Analysis under the weak Morse condition** Now, under the additional Assumptions 4 to 6, it is
 318 interesting to note that a *polynomial dependence* on $1/\epsilon$ is achieved as the following corollary shows.

319 **Corollary 3.2.** *Under the same assumptions as Theorem 3 and Assumptions 4 to 6, taking $\gamma =$
 320 $j(\epsilon) := \frac{4d}{\epsilon} \log\left(\frac{\epsilon L}{M}\right) \vee \frac{8db}{\epsilon^2} \vee 1 \vee \frac{2}{M} \vee C_\gamma$, where C_γ is a constant independent of ϵ defined in Property
 321 C.4, we obtain a gradient complexity of*

$$\tilde{O}\left(\left(n + \frac{n^{\frac{1}{2}}}{\epsilon} \cdot \frac{dL}{C_3} j(\epsilon)\right) C_3^2 j(\epsilon)^4 L^2\right),$$

322 since $\alpha = C_3/\gamma$ (Property C.4).

323 The crux of this corollary is Property C.4. To prove this, we show like Li and Erdogdu (2020) that
 324 ν satisfies the Poincaré inequality with a constant independent of γ . Since it is not hard to show
 325 this around the global minimum, we can step by step extend the set where this inequality holds by a
 326 Lyapunov argument (Theorems D.1 and D.2). The essential difference between this analysis and that
 327 of Li and Erdogdu (2020) is that we do not work on compact manifolds anymore. Some rather minor
 328 difficulties emerge as we cannot employ the compactness but they can be addressed by supposing
 329 dissipativity which assures a quadratic growth for large x .

330 **Remark 2.** *These results do not definitively assert that SARAH-LD and SVRG-LD show the exact
 331 same performance in terms of optimization. Indeed, suppose we are close enough to the global
 332 optimum. Then, a big noise is not necessary anymore since it is more important to stably converge to
 333 the global minimum. Here, we should be able to significantly decrease the noise ϵ_k , and the bottleneck
 334 from the noise should disappear. In this case, SARAH-LD would perform better than SVRG-LD as we
 335 approach the original non-convex optimization setting where SARAH outperforms SVRG.*

336 **Remark 3.** *We also investigated an annealed version of SVRG-LD and SARAH-LD but could not
 337 ameliorate the gradient complexity. The detailed analysis can be found in Appendix E.*

338 5 Discussion and Conclusion

339 The main limitations of our work reside in the gap between practice and theory. Indeed, while our
 340 paper supposes assumptions quite standard in the literature of GLD, it cannot explain the whole
 341 empirical success that machine learning is currently experiencing. Some choices of parameters may
 342 also seem different than the practical use. However, compared to previous work, we succeeded in
 343 the proving convergence of GLD with the popular stochastic gradient with relaxed conditions, and
 344 deleting the dependence of batch size and inner loop length on epsilon, which are all more realistic
 345 situations than prior work. The theoretical study in machine learning and deep learning precisely
 346 plays the role of filling as much as possible this large gap, and our work could be regarded as a further
 347 step forward to achieve this goal. Furthermore, in this paper, we focused on the pure sampling and
 348 optimization performance of the algorithms, and some of the drawbacks are simply due to this fact.
 349 For example, another limitation is that we did not investigate the generalization error in Section 4,
 350 but this was only outside the scope of this work.

351 In conclusion, we analysed the convergence rate of stochastic gradient Langevin Dynamics with
 352 variance reduction under smoothness and LSI and its application to optimization. In Section 3, we
 353 proved the convergence of SVRG-LD in terms of KL-divergence with more relaxed conditions (LSI
 354 and smoothness) and with a better gradient complexity than previous works. We also expanded
 355 SARAH to SARAH-LD and showed that this algorithm enjoyed the same advantages as SVRG-LD
 356 with only an improvement in the coefficients of the gradient complexity. These results led us to
 357 apply SVRG-LD and SARAH-LD to non-convex optimization in Section 4. We provided the global
 358 convergence and a non-asymptotic analysis of SVRG-LD and SARAH-LD. We obtained better
 359 conditions than prior works. Furthermore, we showed that under the additional assumption including
 360 weak Morse and Hessian Lipschitzness, the gradient complexity could be ameliorated, eliminating
 361 the exponential dependence on the inverse of the required error.

362 **References**

- 363 D. Bakry, F. Barthe, P. Cattiaux, and A. Guillin. A simple proof of the Poincaré inequality for a large
364 class of probability measures. *Electronic Communications in Probability*, 13:60–66, 2008.
- 365 K. Balasubramanian, S. Chewi, M. A. Erdogdu, A. Salim, and M. Zhang. Towards a theory of
366 non-log-concave sampling: first-order stationarity guarantees for langevin monte carlo. *arXiv*
367 *preprint arXiv:2202.05214*, 2022.
- 368 V. S. Borkar and S. K. Mitter. A strong approximation theorem for stochastic recursive algorithms.
369 *Journal of Optimization Theory and Applications*, 100(3):499–513, 1999.
- 370 A. Bovier and F. Den Hollander. *Metastability: a Potential-Theoretic Approach*, volume 351. Springer,
371 2016.
- 372 P. Cattiaux, A. Guillin, and L.-M. Wu. A note on Talagrand’s transportation inequality and logarithmic
373 Sobolev inequality. *Probability Theory and Related Fields*, 148(1):285–304, 2010.
- 374 N. Chatterji, N. Flammarion, Y. Ma, P. Bartlett, and M. Jordan. On the theory of variance reduction
375 for stochastic gradient Monte Carlo. In *International Conference on Machine Learning*, pages
376 764–773. PMLR, 2018.
- 377 P. Chen, J. Lu, and L. Xu. Approximation to stochastic variance reduced gradient Langevin dynamics
378 by stochastic delay differential equations. *arXiv preprint arXiv:2106.04357*, 2021.
- 379 T.-S. Chiang, C.-R. Hwang, and S. J. Sheu. Diffusion for global optimization in \mathbb{R}^n . *SIAM Journal*
380 *on Control and Optimization*, 25(3):737–753, 1987.
- 381 A. Dalalyan. Further and stronger analogy between sampling and optimization: Langevin Monte
382 Carlo and gradient descent. In *Conference on Learning Theory*, pages 678–689. PMLR, 2017a.
- 383 A. S. Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave
384 densities. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):
385 651–676, 2017b.
- 386 K. A. Dubey, S. J Reddi, S. A. Williamson, B. Póczos, A. J. Smola, and E. P. Xing. Variance reduction
387 in stochastic gradient Langevin dynamics. *Advances in Neural Information Processing Systems*,
388 29:1154–1162, 2016.
- 389 S. B. Gelfand and S. K. Mitter. Recursive stochastic algorithms for global optimization in \mathbb{R}^d . *SIAM*
390 *Journal on Control and Optimization*, 29(5):999–1018, 1991.
- 391 R. Holley and D. W. Stroock. Logarithmic Sobolev inequalities and stochastic Ising models. 1986.
- 392 Z. Huang and S. Becker. Stochastic gradient Langevin dynamics with variance reduction. In *2021*
393 *International Joint Conference on Neural Networks*, pages 1–8. IEEE, 2021.
- 394 C.-R. Hwang. Laplace’s method revisited: weak convergence of probability measures. *The Annals of*
395 *Probability*, 8(6):1177–1182, 1980.
- 396 R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction.
397 *Advances in Neural Information Processing Systems*, 26:315–323, 2013.
- 398 R. Jordan, D. Kinderlehrer, and F. Otto. The variational formulation of the Fokker–Planck equation.
399 *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998.
- 400 M. B. Li and M. A. Erdogdu. Riemannian langevin algorithm for solving semidefinite programs.
401 *arXiv preprint arXiv:2010.11176v4*, 2020.
- 402 Z. Li. SSRGD: Simple stochastic recursive gradient descent for escaping saddle points. *Advances in*
403 *Neural Information Processing Systems*, 32, 2019.
- 404 J. C. Mattingly, A. M. Stuart, and D. J. Higham. Ergodicity for SDEs and approximations: locally
405 Lipschitz vector fields and degenerate noise. *Stochastic Processes and their Applications*, 101(2):
406 185–232, 2002.

- 407 G. Menz and A. Schlichting. Poincaré and logarithmic Sobolev inequalities by decomposition of the
408 energy landscape. *The Annals of Probability*, 42(5):1809–1884, 2014.
- 409 L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč. SARAH: A novel method for machine learning
410 problems using stochastic recursive gradient. In *International Conference on Machine Learning*,
411 pages 2613–2621. PMLR, 2017a.
- 412 L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč. Stochastic recursive gradient algorithm for
413 nonconvex optimization. *arXiv preprint arXiv:1705.07261*, 2017b.
- 414 F. Otto and C. Villani. Generalization of an inequality by Talagrand and links with the logarithmic
415 Sobolev inequality. *Journal of Functional Analysis*, 173(2):361–400, 2000.
- 416 N. H. Pham, L. M. Nguyen, D. T. Phan, and Q. Tran-Dinh. ProxSARAH: An efficient algorithmic
417 framework for stochastic composite nonconvex optimization. *Journal of Machine Learning
418 Research*, 21(110):1–48, 2020.
- 419 M. Raginsky, A. Rakhlin, and M. Telgarsky. Non-convex learning via stochastic gradient Langevin
420 dynamics: a nonasymptotic analysis. In *Conference on Learning Theory*, pages 1674–1703. PMLR,
421 2017.
- 422 S. Vempala and A. Wibisono. Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry
423 suffices. *Advances in Neural Information Processing Systems*, 32:8094–8106, 2019.
- 424 M. J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, volume 48. Cambridge
425 University Press, 2019.
- 426 Z. Wang, K. Ji, Y. Zhou, Y. Liang, and V. Tarokh. SpiderBoost and momentum: Faster variance
427 reduction algorithms. *Advances in Neural Information Processing Systems*, 32, 2019.
- 428 M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient Langevin dynamics. In
429 *International Conference on Machine Learning*, pages 681–688. Citeseer, 2011.
- 430 A. Wibisono. Sampling as optimization in the space of measures: The Langevin dynamics as a
431 composite optimization problem. In *Conference on Learning Theory*, pages 2093–3027. PMLR,
432 2018.
- 433 P. Xu, J. Chen, D. Zou, and Q. Gu. Global convergence of Langevin dynamics based algorithms for
434 nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- 435 D. Zou, P. Xu, and Q. Gu. Subsampled stochastic variance-reduced gradient Langevin dynamics. In
436 *International Conference on Uncertainty in Artificial Intelligence*, 2018.
- 437 D. Zou, P. Xu, and Q. Gu. Sampling from non-log-concave distributions via variance-reduced gradient
438 Langevin dynamics. In *International Conference on Artificial Intelligence and Statistics*, pages
439 2936–2945. PMLR, 2019a.
- 440 D. Zou, P. Xu, and Q. Gu. Stochastic gradient Hamiltonian Monte Carlo methods with recursive
441 variance reduction. *Advances in Neural Information Processing Systems*, 32:3835–3846, 2019b.
- 442 D. Zou, P. Xu, and Q. Gu. Faster convergence of stochastic gradient Langevin Dynamics for
443 non-log-concave sampling. In *Uncertainty in Artificial Intelligence*, pages 1152–1162. PMLR,
444 2021.

445 Checklist

- 446 1. For all authors...
- 447 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
448 contributions and scope? [Yes]
- 449 (b) Did you describe the limitations of your work? [Yes]
- 450 (c) Did you discuss any potential negative societal impacts of your work? [Yes]

- 451 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
452 them? [Yes]
- 453 2. If you are including theoretical results...
- 454 (a) Did you state the full set of assumptions of all theoretical results? [Yes] All the
455 assumptions are summarized in Subsection 2.3 and referred when used.
- 456 (b) Did you include complete proofs of all theoretical results? [Yes] See the appendices.
- 457 3. If you ran experiments...
- 458 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
459 mental results (either in the supplemental material or as a URL)? [N/A]
- 460 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
461 were chosen)? [N/A]
- 462 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
463 ments multiple times)? [N/A]
- 464 (d) Did you include the total amount of compute and the type of resources used (e.g., type
465 of GPUs, internal cluster, or cloud provider)? [N/A]
- 466 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 467 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 468 (b) Did you mention the license of the assets? [N/A]
- 469 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 470
- 471 (d) Did you discuss whether and how consent was obtained from people whose data you're
472 using/curating? [N/A]
- 473 (e) Did you discuss whether the data you are using/curating contains personally identifiable
474 information or offensive content? [N/A]
- 475 5. If you used crowdsourcing or conducted research with human subjects...
- 476 (a) Did you include the full text of instructions given to participants and screenshots, if
477 applicable? [N/A]
- 478 (b) Did you describe any potential participant risks, with links to Institutional Review
479 Board (IRB) approvals, if applicable? [N/A]
- 480 (c) Did you include the estimated hourly wage paid to participants and the total amount
481 spent on participant compensation? [N/A]