
Mean-Square Analysis with An Application to Optimal Dimension Dependence of Langevin Monte Carlo

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Sampling algorithms based on discretizations of Stochastic Differential Equations
2 (SDEs) compose a rich and popular subset of MCMC methods. This work pro-
3 vides a general framework for the non-asymptotic analysis of sampling error in
4 2-Wasserstein distance, which also leads to a bound of mixing time. The method
5 applies to any consistent discretization of contractive SDEs. When applied to
6 Langevin Monte Carlo algorithm, it establishes $\tilde{O}\left(\sqrt{d}/\epsilon\right)$ mixing time, without
7 warm start, under the common log-smooth and log-strongly-convex conditions,
8 plus a growth condition on the potential of target measures at infinity. This bound
9 improves the best previously known $\tilde{O}\left(d/\epsilon\right)$ result and is optimal in both dimension
10 d and accuracy tolerance ϵ for log-smooth and log-strongly-convex target measures.
11 Our theoretical analysis is further validated by numerical experiments.

12 1 Introduction

13 The problem of sampling statistical distributions has attracted considerable attention, not only in
14 the fields of statistics and scientific computing, but also in machine learning (Robert and Casella,
15 2013; Andrieu et al., 2003; Liu, 2008); for example, how various sampling algorithms scale with
16 the dimension of the target distribution is a popular recent topic in statistical deep learning (see
17 discussions below for references). For samplers that can be viewed as discretizations of SDEs, the
18 idea is to use an ergodic SDE whose equilibrium distribution agrees with the target distribution,
19 and employ an appropriate numerical algorithm that discretizes (the time of) the SDE. The iterates
20 of the numerical algorithm will approximately follow the target distribution when converged, and
21 can be used for various downstream applications such as Bayesian inference and inverse problem
22 (Dashti and Stuart, 2017). One notable example is the Langevin Monte Carlo algorithm (LMC),
23 which corresponds to Euler-Maruyama discretization of overdamped Langevin equation. Its study
24 dated back to at least the 90s (Roberts et al., 1996) but keeps on leading to important discoveries, for
25 example, on non-asymptotics and dimension dependence, which are relevant to machine learning
26 (e.g., Dalalyan (2017a,b); Cheng et al. (2018a); Durmus et al. (2019a,b); Vempala and Wibisono
27 (2019); Dalalyan and Riou-Durand (2020); Erdogdu and Hosseinzadeh (2020); Mou et al. (2019)).
28 LMC is closely related to SGD too (e.g., Mandt et al. (2017)). Many other examples exist, based
29 on alternative SDEs and different discretizations (e.g., Dalalyan and Riou-Durand (2020); Ma et al.
30 (2021); Mou et al. (2021); Li et al. (2020); Roberts and Rosenthal (1998); Chewi et al. (2020); Shen
31 and Lee (2019)).

32 Quantitatively characterizing the non-asymptotic sampling error of numerical algorithms is usually
33 critical for choosing the appropriate algorithm for a specific downstream application, for providing
34 practical guidance on hyperparameter selection and experiment design, and for designing improved
35 samplers. A powerful tool that dates back to (Jordan et al., 1998) is a paradigm of non-asymptotic
36 error analysis, namely to view sampling as optimization in probability space, and it led to many

37 important recent results (e.g., Liu and Wang (2016); Dalalyan (2017a); Wibisono (2018); Zhang
38 et al. (2018); Frogner and Poggio (2020); Chizat and Bach (2018); Chen et al. (2018); Ma et al.
39 (2021); Erdogdu and Hosseinzadeh (2020)). It works by choosing an objective functional, typically
40 some statistical distances/diverges, and showing that the law of the iterates of sampling algorithms
41 converges in that objective functional. However, the choice of the objective functional often needs to
42 be customized for different sampling algorithms. For example, KL divergence works for LMC (Cheng
43 and Bartlett, 2018), but a carefully hand-crafted cross term needs to be added to KL divergence for
44 analyzing KLMC (Ma et al., 2021). Even for the same underlying SDE, different discretization
45 schemes exist and lead to different sampling algorithms, and the analyses of them had usually been
46 case by case (e.g., Cheng et al. (2018b); Dalalyan and Riou-Durand (2020); Shen and Lee (2019)).
47 Therefore, it would be a desirable complement to have a unified, general framework to study the
48 non-asymptotic error of SDE-based sampling algorithms.

49 As an important member of the family of SDE-based sampling algorithms, Langevin Monte Carlo is
50 widely used in practice. Its stochastic gradient version is implemented in common machine learning
51 systems, such as Tensorflow (Abadi et al., 2016), and is the off-the-shelf algorithm for large scale
52 Bayesian inference. With the ever-growing size of parameter space, the non-asymptotic error of LMC
53 is of central theoretical and practical interest, in particular, its dependence on the dimension of the
54 sample space. The best current known upper bound of the mixing time in 2-Wasserstein distance for
55 LMC is $\tilde{O}\left(\frac{d}{\epsilon}\right)$ (Durmus et al., 2019b). Motivated by a recent result (Chewi et al., 2020) that shows
56 better dimension dependence for a Metropolis-Adjusted improvement of LMC, we wonder if the
57 current bound for (unadjusted) LMC is tight, and if not, what is the optimal dimension dependence?

58 **Our contribution** We study a broad family of numerical algorithms that discretize SDEs that
59 have a contraction property (possibly after a coordinate transformation). For this type of problems,
60 we revisit the classical mean-square analysis (Milstein and Tretyakov, 2013) in numerical SDE
61 literature and extend its the global error bound from finite time to infinite time. Same as in classical
62 mean-square analysis, we show the global error is only half order lower than the order of local strong
63 error (p_2). We further obtain a $\tilde{O}\left(C^{\frac{1}{p_2-\frac{1}{2}}}\frac{1}{\epsilon^{\frac{1}{p_2-\frac{1}{2}}}}\right)$ mixing time upper bound in 2-Wasserstein
64 distance for the family of algorithms, where C is a constant containing various information of the
65 underlying problem, e.g., the dimension d .

66 As an application of the general mixing time result, we study the widely used Langevin Monte
67 Carlo algorithm (LMC) for sampling from a Gibbs distribution $\mu \propto \exp(-f(\mathbf{x}))$, which is an
68 Euler-Maruyama discretization of Langevin dynamics. Under the standard smoothness and strong-
69 convexity assumptions, plus an additional linear growth condition on the third-order derivative of f ,
70 we obtain a $\tilde{O}\left(\frac{\sqrt{d}}{\epsilon}\right)$ mixing time in 2-Wasserstein distance, which improves upon the previously best
71 known $\tilde{O}\left(\frac{d}{\epsilon}\right)$ result (Durmus et al., 2019b). For a comparison, note it was known that discretized
72 **kinetic** Langevin dynamics can lead to \sqrt{d} dependence on dimension (Cheng and Bartlett, 2018;
73 Dalalyan and Riou-Durand, 2020) and some believe that it is the introduction of momentum that
74 improves the dimension dependence, but our result shows that discretized overdamped Langevin (no
75 momentum) can also have mixing time scaling like \sqrt{d} . In fact, it is important to mention that it was
76 recently shown that Metropolis-Adjusted Euler-Maruyama discretization of **overdamped** Langevin
77 (i.e., MALA) has an optimal dimension dependence of $\tilde{O}\left(\sqrt{d}\right)$ (Chewi et al., 2020), while what we
78 analyze here is the **unadjusted** version (i.e., LMC), and it has the same dimension dependence (note
79 however that our ϵ dependence is not as good as that for MALA; more discussion in Section 4). We
80 also constructed an example that shows that the mixing time of LMC is at least $\tilde{\Omega}\left(\frac{\sqrt{d}}{\epsilon}\right)$. Hence, our
81 mixing time bound has the optimal dependence on both d and ϵ . Our theoretical analysis is further
82 validated by empirical investigation of numerical examples.

83 2 Preliminaries

84 **Notation** Use the symbol \mathbf{x} to denote a d -dimensional vector, and the plain symbol x to denote a
85 scalar variable. Use $\|\mathbf{x}\|$ to denote the Euclidean norm of vector \mathbf{x} . A numerical algorithm is denoted

86 by \mathcal{A} and its k -th iterate is denoted by $\bar{\mathbf{x}}_k$. We slightly abuse notation by identifying measures with
87 their density function w.r.t. Lebesgue measure. We use the convention $\tilde{\mathcal{O}}(\cdot) = \mathcal{O}(\cdot) \log^{\mathcal{O}(1)}(\cdot)$,
88 i.e., the $\tilde{\mathcal{O}}(\cdot)$ notation ignores the dependence on logarithmic factors. We use the notation $\tilde{\Omega}(\cdot)$
89 similarly. Denote 2-Wasserstein distance by $W_2(\mu_1, \mu_2) = \left(\inf_{(X, Y) \sim \Pi(\mu_1, \mu_2)} \mathbb{E} \|X - Y\|^2 \right)^{\frac{1}{2}}$,
90 where $\Pi(\mu_1, \mu_2)$ is the set of couplings, i.e. all joint measures with X and Y marginals being
91 μ_1 and μ_2 . Denote the target distribution by μ and the law of a random variable \mathbf{X} by $\text{Law}(\mathbf{X})$.
92 Finally, denote the mixing time of an sampling algorithm \mathcal{A} converging to its target distribution μ in
93 2-Wasserstein distance by $\tau_{\text{mix}}(\epsilon; W_2; \mathcal{A}) = \inf \{k \geq 0 \mid W_2(\text{Law}(\bar{\mathbf{x}}_k), \mu) \leq \epsilon\}$.

94 **SDE for Sampling** Consider a general SDE

$$d\mathbf{x}_t = \mathbf{b}(t, \mathbf{x}_t)dt + \boldsymbol{\sigma}(t, \mathbf{x}_t)d\mathbf{B}_t \quad (1)$$

95 where $\mathbf{b} \in \mathbb{R}^d$ is a drift term, $\boldsymbol{\sigma} \in \mathbb{R}^{d \times l}$ is a diffusion coefficient matrix and \mathbf{B}_t is a l -dimensional
96 Wiener process. Under mild condition (Pavliotis, 2014, Theorem 3.1), there exists a unique strong
97 solution \mathbf{x}_t to Eq. (1). Some SDEs admit geometric ergodicity, so that their solutions converge
98 exponentially fast to a unique invariant distribution, and examples include the classical overdamped
99 and kinetic Langevin dynamics, but are not limited to those (e.g., Mou et al. (2021); Li et al. (2020)).
100 Such SDE are desired for sampling purposes, because one can set the target distribution to be the
101 invariant distribution by choosing an SDE with an appropriate potential, and then solve the solution
102 \mathbf{x}_t of the SDE and push the time t to infinity, so that (approximate) samples of the target distribution
103 can be obtained. Except for a few known cases, however, explicit solutions of Eq. (1) are elusive and
104 we have to resort to numerical schemes to simulate/integrate SDE. Such example schemes include,
105 but are not limited to Euler-Maruyama method, Milstein methods and Runge-Kutta method (e.g.,
106 Kloeden and Platen (1992); Milstein and Tretyakov (2013)). With constant stepsize h and at k -th
107 iteration, a typical numerical algorithm takes a previous iterate $\bar{\mathbf{x}}_{k-1}$ and outputs a new iterate $\bar{\mathbf{x}}_k$ as
108 an approximation of the solution \mathbf{x}_t of Eq. (1) at time $t = kh$.

109 **Langevin Monte Carlo Algorithm** LMC algorithm is defined by the following update rule

$$\bar{\mathbf{x}}_k = \bar{\mathbf{x}}_{k-1} - h\nabla f(\bar{\mathbf{x}}_{k-1}) + \sqrt{2h}\boldsymbol{\xi}_k, \quad k = 1, 2, \dots \quad (2)$$

110 where $\{\boldsymbol{\xi}_k\}_{k \in \mathbb{Z}_{>0}}$ are i.i.d. standard d -dimensional Gaussian vectors. LMC corresponds to an Euler-
111 Maruyama discretization of the continuous overdamped Langevin dynamics $d\mathbf{x}_t = -\nabla f(\mathbf{x}_t)dt +$
112 $\sqrt{2}d\mathbf{B}_t$, which converges to an equilibrium distribution $\mu \sim \exp(-f(\mathbf{x}))$.

113 Dalalyan (2017b) provided a non-asymptotic analysis of LMC. An $\tilde{\mathcal{O}}\left(\frac{d}{\epsilon^2}\right)$ mixing time bound in
114 W_2 for log-smooth and log-strongly-convex target measures (Dalalyan, 2017a; Cheng et al., 2018a;
115 Durmus et al., 2019a) has been established. It was further improved to $\tilde{\mathcal{O}}\left(\frac{d}{\epsilon}\right)$ under additional
116 Lipschitz assumption on the Hessian of f (Durmus et al., 2019b). Mixing time bounds of LMC
117 in other statistical distances/divergences have also been studied, including total variation distance
118 (Dalalyan, 2017b; Durmus et al., 2017) and KL divergence (Cheng and Bartlett, 2018).

119 **Classical Mean-Square Analysis** A powerful framework for quantifying the *global* discretization
120 error of a numerical algorithm for Eq. (1), i.e., $e_k = \left\{ \mathbb{E} \|\mathbf{x}_{kh} - \bar{\mathbf{x}}_k\|^2 \right\}^{\frac{1}{2}}$, is mean-square analysis
121 (e.g., the monograph of Milstein and Tretyakov (2013)). Mean-square analysis studies how *local*
122 integration error propagate and accumulate into global integration error; in particular, if one-step
123 (local) weak error and strong error (both the exact solution \mathbf{x}_t and the numerical approximation start
124 from the same initial value \mathbf{x}) satisfy

$$\begin{aligned} \|\mathbb{E}\mathbf{x}_h - \mathbb{E}\bar{\mathbf{x}}_1\| &\leq C_1 \left(1 + \mathbb{E}\|\mathbf{x}\|^2\right)^{\frac{1}{2}} h^{p_1}, & (\text{local weak error}) \\ \left(\mathbb{E}\|\mathbf{x}_h - \bar{\mathbf{x}}_1\|^2\right)^{\frac{1}{2}} &\leq C_2 \left(1 + \mathbb{E}\|\mathbf{x}\|^2\right)^{\frac{1}{2}} h^{p_2}, & (\text{local strong error}) \end{aligned} \quad (3)$$

125 over a time interval $[0, Kh]$ for some constants $C_1, C_2 > 0$, $p_2 \geq \frac{1}{2}$ and $p_1 \geq p_2 + \frac{1}{2}$, then the
126 global error can be bounded by $e_k \leq C \left(1 + \mathbb{E}\|\mathbf{x}_0\|^2\right)^{\frac{1}{2}} h^{p_2 - \frac{1}{2}}$, $k = 1, 2, \dots, K$ for some constant
127 $C > 0$ dependent on Kh .

128 Although classical mean-square analysis is only concerned with numerical integration error, sampling
 129 error can be also inferred. However, there is a limitation that prevents directly employing mean-square
 130 analysis in the non-asymptotic analysis of sampling algorithms. The bound of global error only holds
 131 in finite time because the constant C can grow exponentially as K increases, rendering the bound
 132 useless when $K \rightarrow \infty$.

133 3 Mean-Square Analysis of Samplers Based on Contractive SDE

134 In order to prepare for the analysis of **sampling** error, we first show that the finite time limitation of
 135 **integration** error analysis can be lifted if the SDE being discretized is contractive.

136 More precisely, one bottleneck that prevents the results of classical mean-square analysis from
 137 extending to infinite time horizon, is the fact that the solution of a general SDE may not be bounded,
 138 and neither is its discretization. Note that local error (Eq. (3)) depends on the initial value. To go
 139 from local to global error, these ‘initial’ values correspond to iterates of numerical algorithms, which
 140 change from iteration to iteration and can be unbounded, hence when accumulated together, it is
 141 possible that the global error may blow up.

142 Samplers considered here, on the other hand, are based on stochastic differential equations, each of
 143 which weakly converges to a limiting distributions. The solution of the underlying converging SDE,
 144 as it converges to the invariant measure, gradually inherits boundedness properties from the target
 145 measure. Thus, as long as the target measure has bounded 2nd-moment, a sampling algorithm based
 146 on a reasonable discretization of the SDE should also have bounded 2nd-moment. Motivated by this
 147 observation, we will assume the sampling algorithms we study are based on contractive SDEs, which
 148 is a sufficient condition to ensure the underlying SDE converges to a statistical distribution.

149 **Definition 3.1.** *A stochastic differential equation is contractive if there exists a non-singular constant*
 150 *matrix $A \in \mathbb{R}^{d \times d}$, a constant $\beta > 0$, such that any pair of solutions of the SDE satisfy*

$$\left(\mathbb{E} \|A(\mathbf{x}_t - \mathbf{y}_t)\|^2 \right)^{\frac{1}{2}} \leq \|A(\mathbf{x} - \mathbf{y})\| \exp(-\beta t), \quad (4)$$

151 where $\mathbf{x}_t, \mathbf{y}_t$ are two solutions, driven by the same Brownian motion but evolved respectively from
 152 initial conditions \mathbf{x} and \mathbf{y} .

153 **Remark.** *As long as \mathbf{b} and $\boldsymbol{\sigma}$ in (1) are not explicitly dependent on time, it suffices to find an*
 154 *arbitrarily small $t_0 > 0$ and show (4) holds for all $t < t_0$.*

155 **Remark.** *Sometimes contraction is not easy to establish directly, but can be shown after an appro-*
 156 *priate coordinate transformation, see (Dalalyan and Riou-Durand, 2020, Proposition 1) for such a*
 157 *treatment for kinetic Langevin dynamics. The introduction of A permits such transformations.*

158 We now use contractivity to remove the finite time limitation. We will first need a lemma, which is a
 159 local (short time) result.

160 **Lemma 3.2.** *(Milstein and Tretyakov, 2013, Lemma 1.3) Suppose \mathbf{b} and $\boldsymbol{\sigma}$ in Eq.(1) are Lipschitz*
 161 *continuous. For two solutions $\mathbf{x}_t, \mathbf{y}_t$ of Eq. (1) starting from \mathbf{x}, \mathbf{y} respectively, denote $\mathbf{z} :=$
 162 $(\mathbf{x}_t - \mathbf{x}) - (\mathbf{y}_t - \mathbf{y})$, then there exist $C_0 > 0$ and $h_0 > 0$ such that*

$$\mathbb{E} \|\mathbf{z}\|^2 \leq C_0 \|\mathbf{x} - \mathbf{y}\|^2 t, \quad \forall \mathbf{x}, \mathbf{y}, 0 < t \leq h_0. \quad (5)$$

163 Then we will have a sequence of results that connects **sampling** error (a statistical property) with
 164 local **integration** error (a simulation property). This justifies our generic produce for non-asymptotic
 165 sampling error analysis, which only requires bounding the orders of local weak and strong integration
 166 errors (in addition to establishing contractivity of the continuous dynamics).

167 **Theorem 3.3. (Global Integration Error, Infinite Time Version)** *Suppose Eq.(1) is contractive with*
 168 *rate β and with respect to a non-singular matrix $A \in \mathbb{R}^{d \times d}$, with Lipschitz continuous \mathbf{b} and $\boldsymbol{\sigma}$,*
 169 *and there is a numerical algorithm \mathcal{A} with step size h simulating the solution \mathbf{x}_t of the SDE, whose*
 170 *iterates are denoted by $\bar{\mathbf{x}}_k, k = 0, 1, \dots$. Suppose there exists $0 < h_0 \leq 1, C_1, C_2 > 0, D_1, D_2 \geq$
 171 $0, p_1 \geq 1, \frac{1}{2} < p_2 \leq p_1 - \frac{1}{2}$ such that for any $0 < h \leq h_0$, the algorithm \mathcal{A} has, respectively, local
 172 weak and strong error of order p_1 and p_2 , defined as*

$$\begin{cases} \|\mathbb{E}(\mathbf{x}_h - \bar{\mathbf{x}}_1)\| \leq \left(C_1 + D_1 \sqrt{\mathbb{E} \|\mathbf{x}\|^2} \right) h^{p_1}, \\ \left(\mathbb{E} \|\mathbf{x}_h - \bar{\mathbf{x}}_1\|^2 \right)^{\frac{1}{2}} \leq \left(C_2^2 + D_2^2 \mathbb{E} \|\mathbf{x}\|^2 \right)^{\frac{1}{2}} h^{p_2}, \end{cases} \quad (6)$$

173 where \mathbf{x}_h solves Eq.(1) with any initial value \mathbf{x} and $\bar{\mathbf{x}}_1$ is the result of applying \mathcal{A} to \mathbf{x} for one step.

174 If the solution of SDE \mathbf{x}_t and algorithm \mathcal{A} both start from \mathbf{x}_0 , then for $0 < h \leq h_1 \triangleq$

175 $\min \left\{ h_0, \frac{1}{4\beta}, \left(\frac{\sqrt{\beta}}{4\sqrt{2}\kappa_A D_2} \right)^{\frac{1}{p_2 - \frac{1}{2}}}, \left(\frac{\beta}{8\sqrt{2}\kappa_A (D_1 + C_0 D_2)} \right)^{\frac{1}{p_2 - \frac{1}{2}}} \right\}$, the global error e_k is bounded as

176
$$e_k := \left(\mathbb{E} \|\mathbf{x}_{kh} - \bar{\mathbf{x}}_k\|^2 \right)^{\frac{1}{2}} \leq Ch^{p_2 - \frac{1}{2}}, \quad k = 0, 1, 2, \dots \quad (7)$$

177 where

$$C = \frac{2}{\sqrt{\beta}} \kappa_A^2 \left(\frac{C_1 + C_0 C_2 + \sqrt{2}U(D_1 + C_0 D_2)}{\sqrt{\beta}} + C_2 + \sqrt{2}D_2 U \right), \quad (8)$$

178 C_0 is from Eq. (5), κ_A is the condition number of matrix A and $U^2 \triangleq 4 \|\mathbf{x}_0\|^2 + 5\mathbb{E}_\mu \|\mathbf{x}\|^2$.

179 **Remark.** We use the convention $1/0 = \infty$ when $D_1 = D_2 = 0$. This is pertinent when a numerical
180 algorithm \mathcal{A} , e.g. LMC (Lemma D.3), produces bounded iterates. In such cases, the initial value in
181 Eq. (6) are iterations of \mathcal{A} and will be bounded, it then can be absorbed into C_1, C_2 and we may set
182 $D_1 = D_2 = 0$.

183 Following Theorem 3.3, we obtain the following non-asymptotic bound of the sampling error in W_2 :

184 **Theorem 3.4. (Non-Asymptotic Sampling Error Bound: General Case)** Under the same assump-
185 tion and with the same notation of Theorem 3.3, we have

$$W_2(\text{Law}(\bar{\mathbf{x}}_k), \mu) \leq \sqrt{2}e^{-\beta kh} W_2(\text{Law}(\mathbf{x}_0), \mu) + \sqrt{2}Ch^{p_2 - \frac{1}{2}}, \quad \forall 0 < h \leq h_1.$$

186 A corollary of Theorem 3.4 is a bound on the mixing time of the sampling algorithm:

187 **Corollary 3.5. (Upper Bound of Mixing Time: General Case)** Under the same assumption and
188 with the same notation of Theorem 3.3, we have

$$\tau_{\text{mix}}(\epsilon; W_2; \mathcal{A}) \leq \max \left\{ \frac{1}{\beta h_1}, \frac{1}{\beta} \left(\frac{2C}{\epsilon} \right)^{\frac{1}{p_2 - \frac{1}{2}}} \right\} \log \frac{2\sqrt{2}W_2(\text{Law}(\mathbf{x}_0), \mu)}{\epsilon}$$

189 In particular, when high accuracy is needed, i.e., $\epsilon < 2Ch_1^{p_2 - \frac{1}{2}}$, we have

$$\tau_{\text{mix}}(\epsilon; W_2; \mathcal{A}) \leq \frac{(2C)^{\frac{1}{p_2 - \frac{1}{2}}}}{\beta} \frac{1}{\epsilon^{\frac{1}{p_2 - \frac{1}{2}}}} \log \frac{2\sqrt{2}W_2(\text{Law}(\mathbf{x}_0), \mu)}{\epsilon} = \tilde{\mathcal{O}} \left(\frac{C^{\frac{1}{p_2 - \frac{1}{2}}}}{\beta} \frac{1}{\epsilon^{\frac{1}{p_2 - \frac{1}{2}}}} \right) \quad (9)$$

190 Corollary 3.5 states how mixing time depends on the order of local (strong) error (i.e., p_2) of a
191 numerical algorithm. The larger p_2 is, the shorter the mixing time of the algorithm is, in term of
192 the dependence on accuracy tolerance parameter ϵ . It is important to note that for constant stepsize
193 discretizations that are deterministic on the filtration of the driving Brownian motion and use only its
194 increments, there is a strong order barrier, namely $p_2 \leq 1.5$ (Rüemelin, 1982); however, methods
195 involving multiple stochastic integrals (e.g., Kloeden and Platen (1992); Milstein and Tretyakov
196 (2013)) and randomization (e.g., Shen and Lee (2019)) can yield a larger p_2 .

197 The constant C defined in Eq. (7) typically contains rich information about the underlying SDE, e.g.
198 dimension, Lipschitz constant of drift and noise diffusion, and the initial value \mathbf{x}_0 of the sampling
199 algorithm. Through C , we can uncover the dependence of mixing time bound on various parameters,
200 such as the dimension d . This will be exemplified with Langevin Monte Carlo in the next section.

201 4 Non-Asymptotic Analysis of Langevin Monte Carlo Algorithm

202 This section quantifies how LMC samples from Gibbs target distribution $\mu \sim \exp(-f(\mathbf{x}))$ that has
203 a finite second moment, i.e., $\int_{\mathbb{R}^d} \|\mathbf{x}\|^2 d\mu < \infty$. Assume without loss of generality that the origin is
204 a local minimizer of f , i.e. $\nabla f(\mathbf{0}) = \mathbf{0}$; this is for notational convenience in the analysis and can
205 be realized via a simple coordinate shift, and it is not needed in the practical implementation. In
206 addition, we assume the following two conditions hold:

207 **A 1. (Smoothness and Strong Convexity)** Assume $f \in \mathcal{C}^2$ and is L -smooth and m -strongly-convex,
 208 i.e. there exists $0 < m \leq L$ such that $mI_d \preceq \nabla^2 f(\mathbf{x}) \preceq LI_d$, $\forall \mathbf{x} \in \mathbb{R}^d$.

209 Denote the condition number of f by $\kappa \triangleq \frac{L}{m}$. The smoothness and strong-convexity assumption is
 210 the standard assumption in the literature of analyzing LMC algorithm (Dalalyan, 2017a,b; Cheng and
 211 Bartlett, 2018; Durmus et al., 2019a,b).

212 **A 2. (Linear Growth of the 3rd-order derivative)** Assume $f \in \mathcal{C}^3$ and the operator $\nabla(\Delta f)$ grows
 213 at most linearly, i.e., there exists a constant $G > 0$ such that $\|\nabla(\Delta f(\mathbf{x}))\| \leq G(1 + \|\mathbf{x}\|)$.

214 **Remark.** The linear growth (at infinity) condition on $\nabla\Delta f$ is actually not as restrictive as it appears,
 215 and in some sense even weaker than some classical condition for the existence of solutions to SDE.
 216 For example, a standard condition for ensuring the existence and uniqueness of a global solution to
 217 SDE is at most a linear growth (at infinity) of the drift (Pavliotis, 2014, Theorem 3.1). If we consider
 218 monomial potentials, i.e., $f(x) = x^p$, $p \in \mathbb{N}_+$, then the linear growth condition on $\nabla\Delta f$ is met when
 219 $p \leq 4$, whereas the classical condition for the existence of solutions holds only when $p \leq 2$.

220 To apply mean-square analysis to study LMC algorithm, we will need to ensure the underlying
 221 Langevin dynamics is contractive, which we verify in Section C and D in the appendix. In addition,
 222 we work out all required constants to determine the C in Eq. 7 explicitly in the appendix. With all
 223 these necessary ingredients, we now invoke Theorem 3.4 and obtain the following result:

224 **Theorem 4.1. (Non-Asymptotic Error Bound: LMC)** Suppose Assumption 1 and 2 hold. LMC
 225 iteration $\bar{\mathbf{x}}_{k+1} = \bar{\mathbf{x}}_k - h\nabla f(\bar{\mathbf{x}}_k) + \sqrt{2h}\xi_k$ satisfies

$$W_2(\text{Law}(\bar{\mathbf{x}}_k), \mu) \leq \sqrt{2}e^{-m\kappa h}W_2(\text{Law}(\mathbf{x}_0), \mu) + \sqrt{2}C_{LMC}h, \quad 0 < h \leq \frac{1}{4\kappa L}, k \in \mathbb{N} \quad (10)$$

226 where $C_{LMC} = \frac{10(L^2+G)}{m^{\frac{3}{2}}}\sqrt{2d+m(\|\mathbf{x}_0\|^2+1)} = \mathcal{O}(\sqrt{d})$.

227 Corollary 3.5 combined with the above result gives the following bound on the mixing time of LMC:

228 **Theorem 4.2. (Upper Bound of Mixing Time: LMC)** Suppose Assumption 1 and 2 hold. If running
 229 LMC from \mathbf{x}_0 , we then have

$$\tau_{\text{mix}}(\epsilon; W_2; \text{LMC}) \leq \max\{4\kappa^2, \frac{2C_{LMC}}{m}\frac{1}{\epsilon}\} \log \frac{2\sqrt{2}W_2(\text{Law}(\mathbf{x}_0), \mu)}{\epsilon}$$

230 where C_{LMC} is the same in Theorem 4.1. When high accuracy is needed, i.e., $\epsilon \leq \frac{C_{LMC}}{2m\kappa^2}$, we have

$$\tau_{\text{mix}}(\epsilon; W_2; \text{LMC}) \leq \frac{2C_{LMC}}{m}\frac{1}{\epsilon} \log \frac{2\sqrt{2}W_2(\text{Law}(\mathbf{x}_0), \mu)}{\epsilon} = \tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\epsilon}\right).$$

231 The $\tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\epsilon}\right)$ mixing time bound in 2-Wasserstein distance improves upon the previous ones
 232 (Dalalyan, 2017a; Cheng and Bartlett, 2018; Durmus et al., 2019b,a) in the dependence of d and/or ϵ .
 233 If further assuming $G = \mathcal{O}(L^2)$, we then have $C_{LMC} = \mathcal{O}(\kappa^2\sqrt{m}\sqrt{d})$ and Thm.4.2 shows the mixing
 234 time is $\tilde{\mathcal{O}}\left(\frac{\kappa^2}{\sqrt{m}}\frac{\sqrt{d}}{\epsilon}\right)$, which also improves the κ dependence in some previous results (Dalalyan,
 235 2017a; Cheng and Bartlett, 2018) in the regime $m \leq 1$. A brief comparison is summarized in Table 1.

236 **Optimality** In fact, the $\tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\epsilon}\right)$ mixing time of LMC has the optimal scaling one can expect. This
 237 is in terms of the dependence on d and ϵ , over the class of all log-smooth and log-strongly-convex
 238 target measures. To illustrate this, consider the following Gaussian target distribution whose potential
 239 is

$$f(\mathbf{x}) = \frac{m}{2} \sum_{i=1}^d x_i^2 + \frac{L}{2} \sum_{i=d+1}^{2d} x_i^2, \quad \text{with } m = 1, L \geq 4m. \quad (11)$$

240 We now establish a lower bound on the mixing time of LMC algorithm for this target measure.

241 **Theorem 4.3. (Lower Bound of Mixing Time)** Suppose we run LMC for the target measure defined
 242 in Eq. (11) from $\mathbf{x}_0 = \mathbf{1}_{2d}$, then for any choice of step size $h > 0$ within stability limit, we have

$$\tau_{\text{mix}}(\epsilon; W_2; \text{LMC}) \geq \frac{\sqrt{d}}{8\epsilon} \log \frac{\sqrt{d}}{\epsilon} = \tilde{\Omega}\left(\frac{\sqrt{d}}{\epsilon}\right).$$

Table 1: Comparison of mixing time results in 2-Wasserstein distance of LMC with L -smooth and m -strongly-convex potential. Constant step size is used and accuracy tolerance ϵ is small enough.

	mixing time	Additional Assumption
(Dalalyan, 2017a, Theorem 1)	$\tilde{\mathcal{O}}\left(\frac{\kappa^2}{m} \cdot \frac{d}{\epsilon^2}\right)$	N/A
(Cheng and Bartlett, 2018, Theorem 1)	$\tilde{\mathcal{O}}\left(\frac{\kappa^2}{m} \cdot \frac{d}{\epsilon^2}\right)$	N/A
(Durmus et al., 2019a, Corollary 10)	$\tilde{\mathcal{O}}\left(\frac{\kappa}{m} \cdot \frac{d}{\epsilon^2}\right)$	N/A
(Durmus et al., 2019b, Theorem 8)	$\tilde{\mathcal{O}}\left(\frac{d}{\epsilon}\right)^1$	$\ \nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\ \leq \tilde{L} \ \mathbf{x} - \mathbf{y}\ $
This work (Theorem 4.2)	$\tilde{\mathcal{O}}\left(\frac{\kappa^2}{\sqrt{m}} \cdot \frac{\sqrt{d}}{\epsilon}\right)$	Assumption 2 and $G = \mathcal{O}(L^2)^2$

243 Combining Theorem 4.2 and 4.3, we see that mean-square analysis provides a tight bound for LMC.

244 However, there is one limitation of our result – Assumption 2, which is, although mild, still extra to the
 245 standard setup. Therefore, the gap between the upper bound and the lower bound of LMC algorithm
 246 over the entire family of log-smooth and log-strongly-convex target measures is not completely
 247 closed. We tend to believe that Assumption 2 may not be essential, but rather than an artifact of our
 248 proof technique. We hope to lift this restriction in future work.

249 **Comparison** At least two sampling algorithms are closely related to LMC. One is Kinetic Langevin
 250 Monte Carlo algorithm (KLMC), which is discretized kinetic/underdamped Langevin dynamics, and
 251 the other is Metropolis-Adjusted Langevin Algorithm (MALA) which uses the one-step update of
 252 LMC as a proposal and then accepts/rejects them with a Metropolis-Hastings algorithm.

253 The $\tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\epsilon}\right)$ mixing time in 2-Wasserstein distance of KLMC has been established for log-smooth
 254 and log-strongly-convex target measures in existing literature (Cheng et al., 2018b; Dalalyan and
 255 Riou-Durand, 2020). Due to its better dimension dependence over previously best known results of
 256 LMC, KLMC is understood to be the analog of Nesterov’s accelerated gradient method for sampling
 257 (Ma et al., 2021). Our findings show that LMC is able to achieve the same mixing time, albeit under
 258 an additional growth-at-infinity condition. However, this does not say anything about whether/how
 259 KLMC accelerates LMC, as the optimality of KLMC bound is not yet clear. We also note KLMC has
 260 better condition number dependence, although the κ dependence in our bound may not be tight.

261 In terms of MALA, a recent work (Chewi et al., 2020) establishes a $\tilde{\mathcal{O}}\left(\sqrt{d}\right)$ mixing time in 2-
 262 Wasserstein distance with warm start, and the dimension dependence is shown to be optimal. We see
 263 that without the Metropolis adjustment, LMC can also achieve the optimal dimension dependence as
 264 MALA. But unlike LMC, MALA only has logarithmic dependence on $\frac{1}{\epsilon}$. Under warm-start condition,
 265 is it possible/how to improve the dependence of $\frac{1}{\epsilon}$ for LMC, from polynomial to logarithmic? This
 266 question is beyond the scope of this paper but worth further investigation.

267 5 Numerical Examples

268 This section numerically verifies our theoretical findings for LMC in Section 4, with a particular
 269 focus on the dependence of the discretization error in Theorem 4.1 on dimension d and step size h .
 270 To this end, we consider two target measures specified by the following two potentials:

$$f_1(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2 + \log\left(\sum_{i=1}^d e^{x_i}\right) \quad \text{and} \quad f_2(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2 - \frac{1}{2d^{\frac{1}{2}}} \sum_{i=1}^d \cos\left(d^{\frac{1}{4}} x_i\right). \quad (12)$$

271 It is not hard to see that f_1 is 2-smooth and 1-strongly convex, f_2 is $\frac{3}{2}$ -smooth and 1-strongly-
 272 convex. f_2 is also used in (Chewi et al., 2020) to illustrate the optimal dimension dependence
 273 of MALA. Explicit expression of 2-Wasserstein distance between non-Gaussian distributions is

¹The dependence on κ is not readily available from Theorem 8 in Durmus et al. (2019b).

²The $G = \mathcal{O}(L^2)$ assumption is only for κ, m dependence. Removing it does not affect d, ϵ dependence.

274 typically not available, instead, we use the Euclidean norm of the mean error as a surrogate because
 275 $\|\mathbb{E}\bar{\mathbf{x}}_k - \mathbb{E}_\mu \mathbf{x}\| \leq W_2(\text{Law}(\bar{\mathbf{x}}_k), \mu)$ due to Jensen’s inequality. To obtain an accurate estimate of the
 276 ground truth, we run 10^8 independent LMC realizations using a tiny step size ($h = 0.001$), each till a
 277 fixed, long enough time, and use the empirical average to approximate $\mathbb{E}_\mu \mathbf{x}$.

278 To study the dimension dependence of sampling error, we fix step size $h = 0.1$, and for each
 279 $d \in \{1, 2, 5, 10, 20, 50, 100, 200, 500, 1000\}$, we simulate 10^4 independent Markov chains using
 280 LMC algorithm for 100 iterations, which is long enough for the chain to be well-mixed. The mean
 281 and the standard deviation of the sampling error corresponding to the last 10 iterates are recorded.

282 To study step size dependence of sampling error, dimension is fixed to be $d = 10$. We experiment with
 283 step size $h \in \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \times 10^{-1}$. We fix a continuous time $T = 20$, and run LMC
 284 algorithm for $\lceil \frac{T}{h} \rceil$ iterations for each h . The procedure is repeated 10^4 times with different random
 285 seeds to obtain independent samples. When the corresponding continuous time $t = kh > 10$, we see
 286 from Eq. (10) that LMC is well converged and the sampling error is saturated by the discretization
 287 error. Therefore, for each h , we take the last $\lceil \frac{10}{h} \rceil$ iterates and record the mean and standard deviation
 288 of their sampling error.

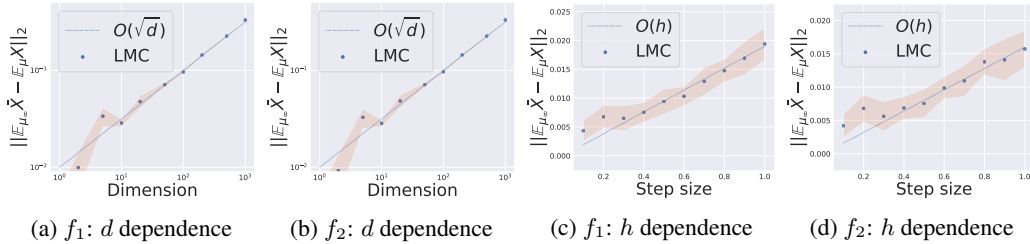


Figure 1: (a) Dependence of the sampling error of LMC on dimension d and step size h for f_1 and f_2 . Both axes in Figure 1a and 1b are in log scale. The shaded areas in Figure 1a and 1b represent one standard deviation of the last 10 iterates. The shaded areas in Figure 1c and 1d represent one standard deviation of the last $\lceil \frac{10}{h} \rceil$ iterations.

289 The experiment results shown in Figure 1 are consistent with our theoretical analysis of the sampling
 290 error. Both linear dependence on \sqrt{d} and h can be identified in and supported by the empirical
 291 evidence. Note results with smaller h are less accurate because one starts to see the error of empirical
 292 approximation due to finite samples. Experiments were conducted on a machine with a 2.20GHz
 293 Intel(R) Xeon(R) E5-2630 v4 CPU and an Nvidia GeForce GTX 1080 GPU.

294 6 Conclusion

295 This paper extends the mean-square analysis framework for analyzing the integration error of SDE
 296 to analyzing the sampling error in 2-Wasserstein distance. Corresponding mixing time bound
 297 unveils how a high-order numerical algorithm can help improve dependence on accuracy tolerance ϵ ,
 298 and potentially other parameters, such as the dimension. When applied to Langevin Monte Carlo
 299 algorithm, it obtains an improved and optimal $\tilde{O}\left(\sqrt{d}/\epsilon\right)$ bound, which was previously thought to be
 300 obtainable only with the addition of momentum.

301 Here are some possible directions worth further investigations. (i) In data-intensive applications,
 302 stochastic gradients are typically used for better scalability. It seems natural to apply the mean-square
 303 analysis framework to study SDE-based stochastic gradient MCMC methods; (ii) Assumption 2 is
 304 likely to be an artifact of our analysis; how to establish the optimal mixing time bound in the standard
 305 log-smooth and log-strongly-convex setup is still an open question; (iii) Motivated by the recent
 306 result of MALA (Chewi et al., 2020), it would be interesting to know whether the dependence on $\frac{1}{\epsilon}$
 307 can be improved to logarithmic, for example if LMC is initialized at a warm start.

308 **References**

- 309 Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G.,
310 Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D. G., Steiner, B., Tucker,
311 P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., and Zheng, X. (2016). Tensorflow: A system
312 for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and
313 Implementation (OSDI 16)*, pages 265–283.
- 314 Andrieu, C., De Freitas, N., Doucet, A., and Jordan, M. I. (2003). An introduction to mcmc for
315 machine learning. *Machine learning*, 50(1):5–43.
- 316 Chen, C., Zhang, R., Wang, W., Li, B., and Chen, L. (2018). A unified particle-optimization frame-
317 work for scalable bayesian sampling. In *The Conference on Uncertainty in Artificial Intelligence*.
- 318 Cheng, X. and Bartlett, P. L. (2018). Convergence of langevin mcmc in kl-divergence. *PMLR 83*,
319 (83):186–211.
- 320 Cheng, X., Chatterji, N. S., Abbasi-Yadkori, Y., Bartlett, P. L., and Jordan, M. I. (2018a). Sharp con-
321 vergence rates for langevin dynamics in the nonconvex setting. *arXiv preprint arXiv:1805.01648*.
- 322 Cheng, X., Chatterji, N. S., Bartlett, P. L., and Jordan, M. I. (2018b). Underdamped langevin mcmc:
323 A non-asymptotic analysis. *Proceedings of the 31st Conference On Learning Theory, PMLR*.
- 324 Chewi, S., Lu, C., Ahn, K., Cheng, X., Gouic, T. L., and Rigollet, P. (2020). Optimal dimension
325 dependence of the metropolis-adjusted langevin algorithm. *arXiv preprint arXiv:2012.12810*.
- 326 Chizat, L. and Bach, F. (2018). On the global convergence of gradient descent for over-parameterized
327 models using optimal transport. In *Advances in neural information processing systems*, pages
328 3036–3046.
- 329 Dalalyan, A. S. (2017a). Further and stronger analogy between sampling and optimization: Langevin
330 monte carlo and gradient descent. *Conference on Learning Theory*, pages 678–689.
- 331 Dalalyan, A. S. (2017b). Theoretical guarantees for approximate sampling from smooth and log-
332 concave densities. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
333 79(3):651–676.
- 334 Dalalyan, A. S. and Riou-Durand, L. (2020). On sampling from a log-concave density using kinetic
335 Langevin diffusions. *Bernoulli*, 26(3):1956–1988.
- 336 Dashti, M. and Stuart, A. M. (2017). *The Bayesian Approach to Inverse Problems*, pages 311–428.
337 Springer International Publishing, Cham.
- 338 Durmus, A., Majewski, S., and Miasojedow, B. (2019a). Analysis of langevin monte carlo via convex
339 optimization. *Journal of Machine Learning Research*, 20:73–1.
- 340 Durmus, A., Moulines, E., et al. (2017). Nonasymptotic convergence analysis for the unadjusted
341 langevin algorithm. *Annals of Applied Probability*, 27(3):1551–1587.
- 342 Durmus, A., Moulines, E., et al. (2019b). High-dimensional bayesian inference via the unadjusted
343 langevin algorithm. *Bernoulli*, 25(4A):2854–2882.
- 344 Erdogdu, M. A. and Hosseinzadeh, R. (2020). On the convergence of langevin monte carlo: The
345 interplay between tail growth and smoothness. *arXiv preprint arXiv:2005.13097*.
- 346 Frogner, C. and Poggio, T. (2020). Approximate inference with wasserstein gradient flows. In
347 *International Conference on Artificial Intelligence and Statistics*.
- 348 Jordan, R., Kinderlehrer, D., and Otto, F. (1998). The variational formulation of the fokker–planck
349 equation. *SIAM journal on mathematical analysis*, 29(1):1–17.
- 350 Kloeden, P. E. and Platen, E. (1992). Numerical solution of stochastic differential equations. Springer.
- 351 Li, R., Zha, H., and Tao, M. (2020). Hessian-free high-resolution nesterov acceleration for sampling.
352 *arXiv preprint arXiv:2006.09230*.

- 353 Liu, J. S. (2008). *Monte Carlo strategies in scientific computing*. Springer Science & Business Media.
- 354 Liu, Q. and Wang, D. (2016). Stein variational gradient descent: A general purpose bayesian inference
355 algorithm. In *Advances in neural information processing systems*, pages 2378–2386.
- 356 Ma, Y.-A., Chatterji, N. S., Cheng, X., Flammarion, N., Bartlett, P. L., and Jordan, M. I. (2021). Is
357 there an analog of Nesterov acceleration for gradient-based MCMC? *Bernoulli*, 27(3):1942 – 1992.
- 358 Mandt, S., Hoffman, M. D., and Blei, D. M. (2017). Stochastic gradient descent as approximate
359 bayesian inference. *Journal of Machine Learning Research*, 18(134):1–35.
- 360 Milstein, G. N. and Tretyakov, M. V. (2013). *Stochastic numerics for mathematical physics*. Springer
361 Science & Business Media.
- 362 Mou, W., Flammarion, N., Wainwright, M. J., and Bartlett, P. L. (2019). Improved bounds for
363 discretization of langevin diffusions: Near-optimal rates without convexity. *arXiv preprint*
364 *arXiv:1907.11331*.
- 365 Mou, W., Ma, Y.-A., Wainwright, M. J., Bartlett, P. L., and Jordan, M. I. (2021). High-order
366 langevin diffusion yields an accelerated mcmc algorithm. *Journal of Machine Learning Research*,
367 22(42):1–41.
- 368 Pavliotis, G. A. (2014). *Stochastic processes and applications: diffusion processes, the Fokker-Planck*
369 *and Langevin equations*, volume 60. Springer.
- 370 Robert, C. and Casella, G. (2013). *Monte Carlo statistical methods*. Springer Science & Business
371 Media.
- 372 Roberts, G. O. and Rosenthal, J. S. (1998). Optimal scaling of discrete approximations to langevin
373 diffusions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(1):255–
374 268.
- 375 Roberts, G. O., Tweedie, R. L., et al. (1996). Exponential convergence of langevin distributions and
376 their discrete approximations. *Bernoulli*, 2(4):341–363.
- 377 Rüemelin, W. (1982). Numerical treatment of stochastic differential equations. *SIAM Journal on*
378 *Numerical Analysis*, 19(3):604–613.
- 379 Shen, R. and Lee, Y. T. (2019). The randomized midpoint method for log-concave sampling. In
380 *Advances in Neural Information Processing Systems*, pages 2098–2109.
- 381 Vempala, S. and Wibisono, A. (2019). Rapid convergence of the unadjusted langevin algorithm:
382 Isoperimetry suffices. In *Advances in Neural Information Processing Systems*, pages 8092–8104.
- 383 Wibisono, A. (2018). Sampling as optimization in the space of measures: The langevin dynamics as
384 a composite optimization problem. In *Conference On Learning Theory*, pages 2093–3027.
- 385 Zhang, R., Chen, C., Li, C., and Carin, L. (2018). Policy optimization as wasserstein gradient flows.
386 In *International Conference on Machine Learning*, pages 5737–5746.

387 **Checklist**

- 388 1. For all authors...
- 389 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
390 contributions and scope? [Yes]
- 391 (b) Did you describe the limitations of your work? [Yes] Limitation is discussed in the
392 paragraph below Table 1.
- 393 (c) Did you discuss any potential negative societal impacts of your work? [N/A] To the
394 best of our knowledge, this work does not have any potential negative societal impact.
- 395 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
396 them? [Yes]
- 397 2. If you are including theoretical results...
- 398 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Assumptions
399 are fully, clearly stated in theorems.
- 400 (b) Did you include complete proofs of all theoretical results? [Yes] Complete proofs of
401 all theoretical results are provided in supplementary materials.
- 402 3. If you ran experiments...
- 403 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
404 mental results (either in the supplemental material or as a URL)? [Yes] All are provided
405 in supplementary materials.
- 406 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
407 were chosen)? [Yes] Our experiments do not need training. We specify all the hyperpa-
408 rameters needed to run the experiments.
- 409 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
410 ments multiple times)? [Yes] See Figure 1.
- 411 (d) Did you include the total amount of compute and the type of resources used (e.g., type
412 of GPUs, internal cluster, or cloud provider)? [Yes] We report the hardware used to run
413 our experiments.
- 414 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 415 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 416 (b) Did you mention the license of the assets? [N/A]
- 417 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 418
- 419 (d) Did you discuss whether and how consent was obtained from people whose data you’re
420 using/curating? [N/A]
- 421 (e) Did you discuss whether the data you are using/curating contains personally identifiable
422 information or offensive content? [N/A]
- 423 5. If you used crowdsourcing or conducted research with human subjects...
- 424 (a) Did you include the full text of instructions given to participants and screenshots, if
425 applicable? [N/A]
- 426 (b) Did you describe any potential participant risks, with links to Institutional Review
427 Board (IRB) approvals, if applicable? [N/A]
- 428 (c) Did you include the estimated hourly wage paid to participants and the total amount
429 spent on participant compensation? [N/A]