# The Mirror Langevin Algorithm Converges with Vanishing Bias

Ruilin Li LIRUILIN 1993 @ GMAIL.COM

Georgia Institute of Technology & Hudson River Trading

Molei Tao Mtao@gatech.edu

Georgia Institute of Technology, School of Mathematics

Santosh S. Vempala VEMPALA@GATECH.EDU

Georgia Institute of Technology, College of Computing

Andre Wibisono Andre.wibisono@yale.edu

Yale University, Department of Computer Science

Editors: Sanjoy Dasgupta and Nika Haghtalab

#### **Abstract**

The technique of modifying the geometry of a problem from Euclidean to Hessian metric has proved to be quite effective in optimization, and has been the subject of study for sampling. The Mirror Langevin Diffusion (MLD) is a sampling analogue of mirror flow in continuous time, and it has nice convergence properties under log-Sobolev or Poincare inequalities relative to the Hessian metric, as shown by Chewi et al. (2020). In discrete time, a simple discretization of MLD is the Mirror Langevin Algorithm (MLA) studied by Zhang et al. (2020), who showed a biased convergence bound with a non-vanishing bias term (does not go to zero as step size goes to zero). This raised the question of whether we need a better analysis or a better discretization to achieve a vanishing bias. Here we study the Mirror Langevin Algorithm and show it indeed has a vanishing bias. We apply mean-square analysis based on Li et al. (2019) and Li et al. (2022) to show the mixing time bound for MLA under the modified self-concordance condition introduced by Zhang et al. (2020).

**Keywords:** Sampling, mirror descent, Langevin dynamics, Wasserstein distance, discretization, mean-square analysis

#### 1. Introduction

Suppose we wish to sample from a probability distribution  $\nu(x) \propto e^{-f(x)}$  supported on a convex set  $\mathcal{X} \subseteq \mathbb{R}^d$  where  $f \colon \mathcal{X} \to \mathbb{R}$  is differentiable. A popular algorithm is the Unadjusted Langevin Algorithm (ULA), which is a basic discretization of the Langevin Dynamics in continuous time:

$$dX_t = -\nabla f(X_t) dt + \sqrt{2} dW_t.$$

Langevin Dynamics has an optimization interpretation as the gradient flow for minimizing relative entropy (KL divergence) with respect to  $\nu$  using the Wasserstein metric  $W_2$  in the space of probability distributions on  $\mathbb{R}^d$ , starting from the seminal work of Jordan et al. (1998); see also Wibisono (2018). In continuous time, Langevin Dynamics has convergence guarantees in various distances, including  $W_2$  distance, KL divergence, or  $\chi^2$ -divergence, under various conditions, such as strong log-concavity, or functional inequalities such as the Log-Sobolev Inequality (LSI) or Poincaré inequality. In discrete time, ULA is a biased discretization of the Langevin Dynamics, and it has an asymptotic bias which scales with the step size. In particular, by setting a small enough step size,

we can obtain a mixing time bound of ULA which has inverse polynomial dependence on the error threshold; see for example Dalalyan (2017); Durmus and Moulines (2017, 2019); Dalalyan and Karagulyan (2019); Vempala and Wibisono (2019); Li et al. (2019); Li et al. (2022).

In many settings, the problem of interest is non-smooth or constrained (e.g., the  $L_1$  ball or a general polytope), and the basic Langevin algorithm does not apply. In optimization, this is handled effectively (and elegantly) by interior-point methods, which use a convex "barrier" function to define a non-Euclidean metric. The resulting metric is given locally by the Hessian of the barrier function. This method results is a convergence bound that scales as  $\sqrt{d}$  for linear and convex optimization.

It is natural to wonder whether such a modification of the geometry could be useful for sampling. Early evidence of this is the Dikin walk, which replaces the ball walk (a discrete-time implementation of constrained Brownian motion) by using an ellipsoid at each step, defined by the Hessian of the logarithmic barrier function. This walk was shown to converge in  $\tilde{O}(md)$  steps for uniformly sampling a polytope with m facets in d dimension (Kannan and Narayanan, 2012). It was recently refined using a weighted barrier function to improve the convergence time to  $\tilde{O}(d^2)$  (Laddha et al., 2020).

A related approach that also originated in optimization is *mirror descent*, which uses a mirror map (the gradient of the barrier function) to change the geometry favorably, and in the context of sampling, can be seen as a generalization of the Langevin algorithm by changing the metric. For constrained sampling, the Mirror Langevin Dynamics was introduced by Zhang et al. (2020) using a mirror map to constrain the domain; see also Hsieh et al. (2018) for a related approach. Mirror Langevin Dynamics is the Langevin dynamics for sampling from  $\nu$  using the Hessian metric generated by the mirror map. In continuous time, Mirror Langevin Dynamics has nice convergence guarantees under an analogous notion of mirror Poincare inequality relative to the Hessian metric, as shown by Chewi et al. (2020); see also Appendix A for a review.

In discrete time, the Mirror Langevin Algorithm (MLA) is a simple discretization of MLD proposed by Zhang et al. (2020), who showed a biased convergence analysis under relative strong convexity and smoothness, but with a non-vanishing bias (does not go to 0 with step size, but remains a constant). Ahn and Chewi (2020) proposed an alternative discretization method which achieves a vanishing bias, but requires an exact simulation of the Brownian motion with changing covariance. Jiang (2021) further showed a convergence analysis of MLA under mirror version of isoperimetry, but still with a non-vanishing bias. These results raised the question of whether we need a better analysis of MLA, or a better discretization of MLD to achieve a vanishing bias.

In this paper, we study the Mirror Langevin Algorithm and show that it indeed has a vanishing bias. The tool we will use is the mean-square analysis framework, proposed by Li et al. (2019) and then refined by Li et al. (2022); the latter version will be used. It will help establish a biased convergence analysis of MLA under relative smoothness, strong convexity and modified self-concordance; these are a subset of the conditions assumed by Zhang et al. (2020). We show that the bias of MLA with step size h scales as  $\sqrt{h}$ ; this leads to a  $\tilde{O}(d/\epsilon^2)$  mixing time bound for MLA (see Theorem 1 and Corollary 2).

# 2. Algorithm and Problem Set-Up

### 2.1. Problem set-up

Suppose we want to sample from a probability distribution  $\nu$  supported on a convex set  $\mathcal{X} \subseteq \mathbb{R}^d$ . We assume  $\nu$  is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}^d$  and has density  $\nu(x) \propto e^{-f(x)}$  for some differentiable  $f \colon \mathcal{X} \to \mathbb{R}$ .

Let  $\phi \colon \mathcal{X} \to \mathbb{R}$  be a twice-differentiable strictly convex function which is of Legendre type (Rockafellar, 1970). This implies  $\nabla \phi(\mathcal{X}) = \mathbb{R}^d$ , and in particular the gradient map  $\nabla \phi \colon \mathcal{X} \to \mathbb{R}^d$  is bijective. We also have  $\nabla^2 \phi(x) \succ 0$  for all  $x \in \mathcal{X}$ . Moreover, we require that  $\|\nabla \phi(x)\| \to \infty$  and  $\nabla^2 \phi(x) \to \infty$  as x approaches the boundary of  $\mathcal{X}$ . Using the Hessian metric  $\nabla^2 \phi$  on  $\mathcal{X}$  will prevent the iterates from leaving the domain  $\mathcal{X}$ . We call  $\nabla \phi \colon \mathcal{X} \to \mathbb{R}^d$  the mirror map and  $\mathcal{Y} = \nabla \phi(\mathcal{X}) = \mathbb{R}^d$  the dual space.

Let  $\phi^* \colon \mathbb{R}^d \to \mathbb{R}$  be the *dual function* of  $\phi$ , defined by  $\phi^*(y) = \sup_{x \in \mathcal{X}} \langle x, y \rangle - \phi(x)$ . Recall  $\nabla \phi^*(y) = \arg \max_{x \in \mathcal{X}} \langle x, y \rangle - \phi(x)$ , and we have  $\nabla \phi^* = (\nabla \phi)^{-1}$ , so  $\nabla \phi(\nabla \phi^*(y)) = y$  for all  $y \in \mathbb{R}^d$ . Furthermore,  $\nabla^2 \phi(x) = \nabla^2 \phi^*(\nabla \phi(x))^{-1}$  for all  $x \in \mathcal{X}$ .

For a vector  $v \in \mathbb{R}^d$ , let  $||v|| = \sqrt{\langle v, v \rangle}$  be the  $\ell_2$ -norm. For a matrix  $A \in \mathbb{R}^{d \times d}$ , let  $||A||_{\mathrm{HS}} = \sqrt{\mathrm{Tr}(AA^\top)}$  be the Hilbert-Schmidt norm.

# 2.2. Mirror Langevin Algorithm

In this paper we study the Mirror Langevin Algorithm:

$$x_{k+1} = \nabla \phi^* \left( \nabla \phi(x_k) - h \nabla f(x_k) + \sqrt{2h} \sqrt{\nabla^2 \phi(x_k)} \, z_k \right) \tag{1}$$

where h>0 is step size and  $z_k\sim\mathcal{N}(0,I)$  is an independent Gaussian random variable in  $\mathbb{R}^d$ . Here  $\sqrt{\nabla^2\phi(x)}$  is a square-root of  $\nabla^2\phi(x)$ , i.e. any matrix  $C(x)\in\mathbb{R}^{d\times d}$  satisfying  $C(x)C(x)^\top=\nabla^2\phi(x)$ . This algorithm can be seen as a sampling version of the mirror descent algorithm from optimization, since we can write the update of MLA in the following form which resembles mirror descent:

$$x_{k+1} = \arg\min_{x \in \mathcal{X}} \left\{ \langle h \nabla f(x_k) - \sqrt{2h} \sqrt{\nabla^2 \phi(x_k)} z_k, x - x_k \rangle + D_{\phi}(x, x_k) \right\}$$

where  $D_{\phi}(x,x') = \phi(x) - \phi(x') - \langle \nabla \phi(x'), x - x' \rangle$  is the Bregman divergence of  $\phi$ . In particular, in the *Euclidean case*, i.e. when  $\mathcal{X} = \mathbb{R}^d$  and  $\phi(x) = \frac{1}{2} ||x||^2$ , MLA recovers the usual Unadjusted Langevin Algorithm (ULA).

MLA can be seen as a coordinate-transformed Euler-Maruyama discretization of the **Mirror Langevin Dynamics** in continuous time, given by

$$\begin{cases} Y_t &= \nabla \phi(X_t) \\ dY_t &= -\nabla f(X_t) dt + \sqrt{2} \sqrt{\nabla^2 \phi(X_t)} dW_t \end{cases}.$$

See Section 2.3 for a reformulation purely in the dual space, and Appendix A for more details on the continuous-time dynamics. Zhang et al. (2020) studied MLA as a simple discretization of the Mirror Langevin Dynamics, and showed that under certain assumptions, the iterates of MLA converge to a Wasserstein ball around the target with some radius which depends on the modified self-concordance parameter of  $\phi$  (see Section 3.1 for more detail). In the Euclidean case (when  $\phi$ 

is quadratic) this radius is 0, so MLA converges to  $\nu$  with sufficiently small step size, recovering the typical bound for ULA. However, for general  $\phi$ , the radius is positive. Therefore, the result of Zhang et al. (2020) only guarantees MLA enters a ball around the target, but it may not converge to the target even when the step size goes to 0. They further conjectured the bias is unavoidable. This raises an interesting question of whether the non-vanishing bias of MLA is indeed unavoidable because we are discretizing a diffusion process with changing covariance, or whether there is a better analysis of MLA with vanishing bias. Here we show that indeed MLA has a vanishing bias, by applying the mean-square analysis framework of Li et al. (2019) and Li et al. (2022).

# 2.3. Mirror Langevin Algorithm in the Dual Space

Let us work in the dual space  $\mathcal{Y} = \nabla \phi(\mathcal{X}) = \mathbb{R}^d$  via the mirror map  $\nabla \phi \colon \mathcal{X} \to \mathbb{R}^d$ . Given  $x \in \mathcal{X}$ , we define the *dual variable* 

$$y = \nabla \phi(x) \in \mathbb{R}^d$$

and its inverse is given by  $x = \nabla \phi^*(y)$ . The target distribution  $\tilde{\nu}$  on the dual space is the pushforward of the original target  $\nu \propto e^{-f}$  under the mirror map:  $\tilde{\nu} = (\nabla \phi)_{\#} \nu$ . If we write the density as  $\tilde{\nu}(y) \propto e^{-\tilde{f}(y)}$ , then we have  $\tilde{f}(y) = f(\nabla \phi^*(y)) - \log \det \nabla^2 \phi^*(y)$ . Moreover, the Hessian metric  $\nabla^2 \phi(x)$  on  $\mathcal{X}$  corresponds to the Hessian metric  $\nabla^2 \phi^*(y)$  on  $\mathbb{R}^d$  generated by the dual function  $\phi^*$ ; that is,  $\nabla^2 \phi^*$  on  $\mathbb{R}^d$  is the pullback metric of  $\nabla^2 \phi$  on  $\mathcal{X}$  under the inverse mirror map  $\nabla \phi^* \colon \mathbb{R}^d \to \mathcal{X}$ . Therefore, the metric space  $(\mathcal{X}, \nabla^2 \phi)$  is isometric to  $(\mathbb{R}^d, \nabla^2 \phi^*)$ .

If  $x_k \in \mathcal{X}$  follows the Mirror Langevin Algorithm (1), then  $y_k = \nabla \phi(x_k) \in \mathbb{R}^d$  follows the Mirror Langevin Algorithm in the dual space:

$$y_{k+1} = y_k - h\nabla f(\nabla \phi^*(y_k)) + \sqrt{2h}\sqrt{\nabla^2 \phi^*(y_k)^{-1}} z_k.$$
 (2)

MLA in the dual space (2) can be seen as a discretization of the mirror Langevin dynamics to sample from  $\tilde{\nu} \propto e^{-\tilde{f}}$  with the Hessian metric  $\nabla^2 \phi^*$  on  $\mathbb{R}^d$ .

Let us define  $q: \mathbb{R}^d \to \mathbb{R}^d$  and  $A: \mathbb{R}^d \to \mathbb{R}^{d \times d}$  by

$$g(y) = \nabla f(\nabla \phi^*(y)) \tag{3}$$

$$A(y) = \sqrt{\nabla^2 \phi^*(y)^{-1}}. (4)$$

Note here A(y) is any square-root of  $\nabla^2 \phi^*(y)^{-1}$ . Then we can write MLA in the dual space as

$$y_{k+1} = y_k - hg(y_k) + \sqrt{2h}A(y_k)z_k.$$
 (5)

As  $h \to 0$ , MLA converges to the **Mirror Langevin Dynamics**, which is a continuous-time stochastic process  $Y_t \in \mathbb{R}^d$  following the stochastic differential equation:

$$dY_t = -g(Y_t)dt + \sqrt{2}A(Y_t)dW_t$$

where  $W_t$  is the standard Brownian motion in  $\mathbb{R}^d$ ; see Section 4.2.1 for more properties.

# 2.4. Wasserstein distance in dual space

Along MLA in the dual space (2), let  $\tilde{\rho}_k$  denote the distribution of the random variable  $y_k \in \mathbb{R}^d$ . We will show a convergence analysis of MLA in the dual space in terms of the Euclidean Wasserstein distance  $W_2$  between  $\tilde{\rho}_k$  and  $\tilde{\nu}$  on  $\mathbb{R}^d$ :

$$W_2(\tilde{\rho}, \tilde{\nu})^2 = \inf_{y \sim \tilde{\rho}, y^* \sim \tilde{\nu}} \mathbb{E}[\|y - y^*\|^2].$$

Note that this distance does not use the Hessian metric  $\nabla^2 \phi^*$  on  $\mathbb{R}^d$ . In the original space  $\mathcal{X}$ , this gives a modified  $W_2$  distance under the mirror map:

$$W_{2,\phi}(\rho,\nu)^2 = \inf_{x \sim \rho, x' \sim \nu} \mathbb{E}[\|\nabla \phi(x) - \nabla \phi(x')\|^2].$$

That is, if  $\tilde{\rho} = (\nabla \phi)_{\#} \rho$  and  $\tilde{\nu} = (\nabla \phi)_{\#} \nu$ , then  $W_{2,\phi}(\rho,\nu) = W_2(\tilde{\rho},\tilde{\nu})$ . This is the same modified Wasserstein distance that is used in Zhang et al. (2020). This corresponds to using the *squared* Hessian metric  $(\nabla^2 \phi(x))^2$  on  $\mathcal{X}$ , which is isometric to the Euclidean metric I on  $\mathbb{R}^d$  (rather than the Hessian metric  $\nabla^2 \phi(x)$  on  $\mathcal{X}$ , which is isometric to the Hessian metric  $\nabla^2 \phi^*(y)$  on  $\mathbb{R}^d$ , and which is used in the continuous-time analysis in Chewi et al. (2020)).

# 3. Main Result: Mixing Time Bound for MLA

We present our main result on the mixing time bound of MLA. We need the following assumptions.

(A1)  $\phi$  satisfies the **modified self-concordance** property with parameter  $\alpha > 0$ , which means:

$$\|\sqrt{\nabla^2 \phi(x')} - \sqrt{\nabla^2 \phi(x)}\|_{\mathrm{HS}} \le \sqrt{\alpha} \|\nabla \phi(x') - \nabla \phi(x)\|_2 \quad \forall x', x \in \mathcal{X}.$$

Equivalently,  $A(y) = \sqrt{\nabla^2 \phi^*(y)^{-1}}$  is  $\sqrt{\alpha}$ -Lipschitz in the Hilbert-Schmidt norm:

$$||A(y') - A(y)||_{HS} \le \sqrt{\alpha} ||y' - y||_2 \quad \forall y', y \in \mathbb{R}^d.$$

(A2) f is M-smooth with respect to  $\phi$  for some  $0 < M < \infty$ , which means:

$$\|\nabla f(x') - \nabla f(x)\|_2 \le M \|\nabla \phi(x') - \nabla \phi(x)\|_2 \quad \forall x', x \in \mathcal{X}.$$

Equivalently,  $g(y) = \nabla f(\nabla \phi^*(y))$  is M-Lipschitz:

$$||g(y') - g(y)||_2 \le M||y' - y||_2 \quad \forall y', y \in \mathcal{Y}.$$

(A3) f is m-strongly convex with respect to  $\phi$  for some  $0 < m \le M$ , which means:

$$\langle \nabla f(x') - \nabla f(x), \nabla \phi(x') - \nabla \phi(x) \rangle \ge m \|\nabla \phi(x') - \nabla \phi(x)\|_2^2 \quad \forall x', x \in \mathcal{X}.$$

Equivalently,  $g(y) = \nabla f(\nabla \phi^*(y))$  is m-monotone:

$$\langle g(y') - g(y), y' - y \rangle \ge m \|y' - y\|_2^2 \quad \forall y', y \in \mathbb{R}^d.$$

These are a subset of the assumptions in Zhang et al. (2020). In particular, we do not assume a bound on the commutator of  $\nabla^2 f$  and  $\nabla^2 \phi$ . Our main result is the following.

**Theorem 1** Assume (A1), (A2), (A3), and assume  $\alpha < m$ . There is a maximum step size  $h_{\max} = \mathcal{O}\left(\frac{(m-\alpha)^2}{M^2(1+4\alpha)^2}\right)$  and constant  $C_{\text{MLA}} = \mathcal{O}\left(\frac{M(1+4\alpha)\sqrt{d}}{m-\alpha}\right)$ , such that if we run MLA (1) with  $0 < h \le h_{\max}$  from any  $x_0 \sim \rho_0$ , then the iterates  $x_k \sim \rho_k$  satisfy:

$$W_{2,\phi}(\rho_k,\nu) \le e^{-(m-\alpha)hk} W_{2,\phi}(\rho_0,\nu) + C_{\text{MLA}} \sqrt{h}.$$

Equivalently, if we run MLA in the dual space (2) with  $0 < h \le h_1$  from any  $y_0 \sim \tilde{\rho}_0$ , then the iterates  $y_k \sim \tilde{\rho}_k$  satisfy:

$$W_2(\tilde{\rho}_k, \tilde{\nu}) \le e^{-(m-\alpha)hk} W_2(\tilde{\rho}_0, \tilde{\nu}) + C_{\text{MLA}} \sqrt{h}.$$

See Section 4.3 for the proof of Theorem 1 and explicit forms of the constant  $C_{\rm MLA}$  and maximum step size  $h_{\rm max}$ . This result shows MLA has a bias that is vanishing with step size, and thus we can reach an arbitrary accuracy by using a small enough step size. In particular, this improves on the analysis in Zhang et al. (2020), which has a non-vanishing bias and under stronger assumptions. By choosing a small step size, we obtain the following mixing time bound for MLA.

**Corollary 2** For any (small) error threshold  $\epsilon > 0$ , to reach  $W_2(\tilde{\rho}_k, \tilde{\nu}) \leq \epsilon$ , it suffices to run MLA in the dual space (2) with step size  $h = \frac{\epsilon^2}{4C_{\rm MLA}^2}$  for  $k = \tau_{W_2}(\epsilon)$  iterations where

$$\tau_{W_2}(\epsilon) \le \frac{1}{(m-\alpha)h} \log \frac{2W_2(\tilde{\rho}_0, \tilde{\nu})}{\epsilon} = \tilde{O}\left(\frac{C_{\text{MLA}}^2}{(m-\alpha)\epsilon^2}\right) = \tilde{O}\left(\frac{M^2(1+4\alpha)^2d}{(m-\alpha)^3\epsilon^2}\right).$$

#### 3.1. Discussion of result

Theorem 1 shows that MLA has a biased convergence guarantee where the bias scales as  $O(\sqrt{dh})$  where d is dimension and h is step size (assuming  $m, M, \alpha$  are independent of d for now). This leads to a mixing time bound of  $\tilde{O}(d/\epsilon^2)$  for MLA.

Let us compare MLA with ULA (i.e., MLA in the Euclidean case with  $\phi(x) = \frac{1}{2}||x||^2$ ). Recall for ULA, the mean-square analysis by Li et al. (2022) yields a biased convergence guarantee where the bias scales as  $O(\sqrt{dh})$  under an additional 3rd-order regularity condition on f. This leads to a mixing time bound of  $O(\sqrt{d/\epsilon})$  for ULA. We see the bias of MLA has a worse dependence on h than the bias of ULA. This is because the continuous-time Mirror Langevin Dynamics (7) of MLA has a changing covariance, while the usual continuous-time Langevin Dynamics of ULA has a constant covariance; therefore, MLA incurs an additional stochastic error from the Brownian motion part, which is not incurred by ULA. Formally, this is reflected in the orders of error of the two algorithms: We show below that MLA has local weak and strong errors of orders  $p_1 = \frac{3}{2}$  at least and  $p_2 = 1$  (note the local weak order of MLA is actually  $p_1 = 2$ , because it is the Euler-Maruyama discretization of an SDE; the multiplicative noise causes the strong error to lose half an order, but not the weak error (see e.g., (Milstein and Tretyakov, 2013, page 14)); however, we will see that as long as  $p_1 \ge p_2 + \frac{1}{2}$ , the order of the final sampling error is determined by  $p_2$  but not  $p_1$ , and even though our  $p_1 = \frac{3}{2}$  bound is not tight in order, its constants can be made very explicit and hence helpful to later analysis). On the other hand, it is well known that ULA has local weak and strong error of orders  $p_1=2$  and  $p_2=\frac{3}{2}$  because it is the Euler-Maruyama discretization of an SDE with additive noise (see Milstein and Tretyakov (2013) for the general theory and Li et al. (2022) for details of worked out constants). It would be interesting to understand whether we can improve the local errors and the bias of MLA, perhaps using more sophisticated discretization of MLD to improve the stochastic error.

Our result improves on the analysis of Zhang et al. (2020), who assume stronger assumptions (our assumptions (A1), (A2), (A3), along with two assumptions on the moment of  $\nabla^2 \phi$  and a bound on the commutator of  $\nabla^2 f$  and  $\nabla^2 \phi$ ), and prove a biased convergence analysis where the bias scales as  $O(\sqrt{dh} + r_0)$ , where  $r_0 = O(\sqrt{\alpha d})$  does not depend on h. Note in the Euclidean case (when  $\phi(x) = \frac{1}{2}||x||^2$ ), the modified self-concordance parameter is  $\alpha = 0$ , and thus  $r_0 = 0$ ; but for general  $\phi$ , the asymptotic radius is positive:  $r_0 > 0$ , so the result of Zhang et al. (2020) does not guarantee convergence to the target. With our mean-square analysis, we have shown that in fact there is no dependence on this radius  $r_0$ , and the bias indeed scales as  $O(\sqrt{dh})$ .

We note our result uses the modified self-concordance property, as also used in Zhang et al. (2020). In one-dimension (d=1), modified self-concordance is equivalent to the classical self-concordance property: Both are equivalent to the condition that  $x\mapsto 1/\sqrt{\phi''(x)}$  is a Lipschitz function. However, in higher dimension, they are different. In particular, modified self-concordance is not an affine-invariant property (in contrast to the classical self-concordance), and the parameter  $\alpha$  can be arbitrarily large; see example in Appendix D. This is problematic since our convergence bound only holds when  $\alpha$  is less than m (the strong convexity parameter). It would be desirable to have an analysis of MLA with the more natural self-concordance property.

Our result in Theorem 1 shows that to obtain a consistent algorithm (with a vanishing bias) from MLD, it suffices to apply a simple discretization such as MLA. This shows we do not need to use an exact simulator of the Brownian motion with changing covariance, as proposed by Ahn and Chewi (2020), which allows a nice analysis under self-concordance property. It would be interesting to bridge the analysis technique to MLA.

The relative smoothness (A2) and relative strong convexity (A3) conditions imply that the Hessian of f are bounded by the Hessian of  $\phi$ :

$$m\nabla^2\phi(x)\preceq\nabla^2f(x)\preceq M\nabla^2\phi(x) ~~\forall~x\in\mathcal{X}.$$

See (Zhang et al., 2020, Appendix B) for more details. Since we assume  $\phi$  is a Legendre function,  $\nabla^2\phi(x)\to\infty$  as  $x\to\partial\mathcal{X}$ ; then for our result to hold, we need  $\nabla^2f\to\infty$  as  $x\to\partial\mathcal{X}$ . This restricts the applicability of the result; for example, it does not apply when  $\nu$  is a uniform (f=0) or Gaussian distribution (f) is quadratic) restricted on a polytope with  $\phi$  being the log-barrier function. It is desirable to have a more general convergence analysis of MLA under weaker conditions on f and  $\phi$ .

# 4. Proof of main result

The proof of Theorem 1 uses the mean-square analysis framework described in Li et al. (2022). We review the mean-square analysis framework in Section 4.1. We verify the conditions hold for MLA in Section 4.2, and apply the mean-square analysis to prove Theorem 1 in Section 4.3.

#### 4.1. A review of the mean-square analysis framework

Mean-square analysis was a classical tool for analyzing the integration error of SDEs (e.g., Milstein and Tretyakov (2013)). Li et al. (2019) extended it to obtain non-asymptotic sampling error bound of an algorithm which is a discretization of a decaying stochastic differential equation (SDE).

While Li et al. (2019) required the local errors to satisfy uniform bounds, Li et al. (2022) relaxes this requirement and only needs non-uniform bounds. We will establish non-uniform local error bounds for MLA, and thus use the version of mean-square analysis in Li et al. (2022). The results will be reviewed in a simplified setting; see (Li et al., 2022, Section 3) for details.

**Contractive SDE.** Consider a continuous-time process  $Y_t \in \mathbb{R}^d$  which evolves following the SDE:

$$dY_t = -g(Y_t) dt + \sqrt{2}A(Y_t) dW_t \tag{6}$$

for some vector field  $g: \mathbb{R}^d \to \mathbb{R}^d$  and matrix  $A: \mathbb{R}^d \to \mathbb{R}^{d \times d}$ . We assume g and A are Lipschitz continuous. Here  $W_t$  is the standard Brownian motion in  $\mathbb{R}^d$ .

We say the SDE (6) is **contractive** with rate  $\beta > 0$  if there exists  $t_0 > 0$  such that any two solutions  $Y_t, Y_t'$  with synchronous coupling (i.e. driven by the same Brownian motion) satisfy:

$$\mathbb{E}[\|Y_t - Y_t'\|^2] \le e^{-2\beta t} \mathbb{E}[\|Y_0 - Y_0'\|^2] \quad \forall t \in (0, t_0).$$

If the SDE (6) is contractive, then it has a stationary distribution  $\tilde{\nu}$ .

**Short-time deviation.** Since g and A are Lipschitz continuous, one can show (Milstein and Tretyakov, 2013, Lemma 1.3) that there exist a maximum time  $t_0 > 0$  and a constant  $C_0 > 0$  such that for any solutions  $Y_t, Y_t'$  with synchronous coupling:

$$\mathbb{E}[\|(Y_t' - Y_0') - (Y_t - Y_0)\|_2^2] \le C_0 \,\mathbb{E}[\|Y_0' - Y_0\|_2^2] \,t \qquad \forall \, 0 < t \le t_0.$$

**Algorithm and local error.** Suppose we have an algorithm  $Alg_h$  depending on a step size h > 0 that simulates the solution  $Y_t$  of the SDE (6) at time t = h.

For any  $Y_0 \in \mathbb{R}^d$ , let  $Y_h$  denote the solution of the SDE (6) at time t = h, and let  $\bar{Y}_1 = \mathsf{Alg}_h(Y_0)$  denote the output of the algorithm from  $Y_0$ . We say that the algorithm has (non-uniform) **local weak error** of order  $p_1$  if there exist a maximum step size  $h_1 > 0$  and constants  $C_1, D_1 \ge 0$  such that

$$\|\mathbb{E}[Y_h - \bar{Y}_1]\| \le \left(C_1 + D_1 \sqrt{\mathbb{E}[\|Y_0\|^2}\right) h^{p_1} \quad \forall \ 0 < h \le h_1.$$

We say the algorithm has (non-uniform) **local strong error** of order  $p_2$  if there exist a maximum step size  $h_2 > 0$  and constants  $C_2, D_2 \ge 0$  such that

$$\mathbb{E}[\|Y_h - \bar{Y}_1\|^2] \le \left(C_2^2 + D_2^2 \mathbb{E}[\|Y_0\|^2\right) h^{2p_2} \quad \forall \ 0 < h \le h_2.$$

Here  $Y_h$  and  $\bar{Y}_1$  are coupled by sharing the same filtration (i.e. the algorithm  $\mathsf{Alg}_h$  has access to the realization of the Wiener process that generates  $Y_h$ ).

When  $D_1 = D_2 = 0$ , the bounds are termed as uniform bounds in Li et al. (2019).

**Bound on global error.** With the set-up above, the mean-square analysis framework produces the following bound on the global (long-term) error.

**Theorem 3** ((Li et al., 2022, Theorem 3.3, 3.4)) Assume the SDE (6) is contractive with rate  $\beta > 0$ . Assume the algorithm  $Alg_h$  has local weak error of order  $p_1$  and local strong error of order  $p_2$  with  $\frac{1}{2} < p_2 \le p_1 - \frac{1}{2}$ . Let us define a maximum step size  $h_{max} > 0$  by

$$h_{\max} = \min \left\{ t_0, h_1, h_2, \frac{1}{4\beta}, \left( \frac{\sqrt{\beta}}{4\sqrt{2}D_2} \right)^{\frac{1}{p_2 - \frac{1}{2}}}, \left( \frac{\beta}{8\sqrt{2}(D_1 + C_0 D_2)} \right)^{\frac{1}{p_2 - \frac{1}{2}}} \right\}$$

and constants  $U = \sqrt{4\mathbb{E}[\|Y_0\|^2] + 6\mathbb{E}_{\tilde{\nu}}[\|Y\|^2]}$  and C > 0 by

$$C = \frac{2}{\sqrt{\beta}} \left( \frac{C_1 + C_0 C_2 + \sqrt{2}U(D_1 + C_0 D_2)}{\sqrt{\beta}} + C_2 + \sqrt{2}D_2 U \right).$$

Starting from any  $Y_0 = \bar{Y}_0 \sim \tilde{\rho}_0$ , suppose we run the algorithm  $\mathsf{Alg}_h$  with step size  $0 < h \le h_{\max}$  to produce iterates  $\bar{Y}_k = \mathsf{Alg}_h(\bar{Y}_{k-1}) \sim \tilde{\rho}_k$ . Let  $Y_{hk}$  denote the solution to the SDE (6) at time t = hk. Then  $\bar{Y}_k$  is close to  $Y_{hk}$  at all time:

$$\sqrt{\mathbb{E}[\|Y_{hk} - \bar{Y}_k\|^2]} \le Ch^{p_2 - \frac{1}{2}} \quad \forall k \ge 0.$$

Furthermore, the distribution of  $\bar{Y}_k \sim \tilde{\rho}_k$  has the following biased convergence guarantee:

$$W_2(\tilde{\rho}_k, \tilde{\nu}) \le e^{-\beta kh} W_2(\tilde{\rho}_0, \tilde{\nu}) + Ch^{p_2 - \frac{1}{2}} \quad \forall k \ge 0.$$

### 4.2. Application to MLA

For our sampling problem, we wish to apply the mean-square analysis framework to the Mirror Langevin Algorithm in the dual space (2). The continuous-time SDE (6) of MLA is the Mirror Langevin Dynamics, which we review in the next section. We establish the local error orders of MLA in the following section.

#### 4.2.1. MIRROR LANGEVIN DYNAMICS

Consider the Mirror Langevin Dynamics (MLD), which is a stochastic process  $Y_t \in \mathbb{R}^d$  following the SDE:

$$dY_t = -g(Y_t)dt + \sqrt{2}A(Y_t)dW_t \tag{7}$$

where as defined in (3) and (4),  $g(y) = \nabla f(\nabla \phi^*(y))$  and  $A(y) = \sqrt{\nabla^2 \phi^*(y)^{-1}}$ . The stationary distribution of MLD (7) is the target distribution in the dual space:  $\tilde{\nu} = (\nabla \phi)_{\#} \nu$ .

By assumptions (A1) and (A2), g and A are Lipschitz continuous. Let us establish the contractivity and deviation bound on MLD. The proofs are provided in Appendix B.

**Lemma 4** Assume (A1) and (A2) with  $\alpha < m$ . Then MLD (7) is contractive with rate  $\beta = m - \alpha$ .

**Lemma 5** Assume (A1), (A2), and (A3) with  $\alpha < m$ . Then any two solutions  $Y_t, Y'_t$  of MLD (7) with synchronous coupling satisfy

$$\mathbb{E}[\|(Y_t' - Y_0') - (Y_t - Y_0)\|^2] \le 4M \,\mathbb{E}[\|Y_0' - Y_0\|^2] \,t \qquad \forall \, t \ge 0.$$

We also need the following bound on MLD. Let  $x^* = \arg\min_{x \in \mathcal{X}} f(x)$  and  $y^* = \nabla \phi(x^*) \in \mathbb{R}^d$ .

**Lemma 6** Assume (A1), (A2), and (A3). Along MLD (7), for  $0 < t \le \frac{1}{M^2 + 4\alpha}$ ,

$$\mathbb{E}[\|Y_t - Y_0\|^2] \le \gamma t \tag{8}$$

where 
$$\gamma = 8(1+4\alpha)\mathbb{E}[\|Y_0\|^2] + 8(1+4\alpha)\|y^*\|^2 + 16\|A(y^*)\|_{HS}^2 + \frac{4}{M^2}\|g(y^*)\|^2$$
.

**Remark 7** In Lemma 4 we show MLD is contracting if  $\alpha < m$ . In general, a bound on  $\alpha$  (the Lipschitz constant of the covariance) is necessary for an SDE with multiplicative noise to contract; see the example of the geometric Brownian motion in Appendix C.

#### 4.2.2. LOCAL ERRORS OF THE MIRROR LANGEVIN ALGORITHM

Let us now consider the algorithm  $Alg_h$  to be the Mirror Langevin Algorithm in the dual space (2). We can show MLA has the following local errors. The proofs are provided in Appendix B.

**Lemma 8** Assume (A1), (A2), and (A3). Then MLA (2) has local weak error at least of order  $p_1 = \frac{3}{2}$ , with maximum step size  $h_1 = \frac{1}{M^2 + 4\alpha}$  and constants

$$C_1 = 3M\sqrt{1+4\alpha} \left( \|y^*\| + \|A(y^*)\|_{HS} + \frac{1}{M} \|g(y^*)\| \right)$$
$$D_1 = 2M\sqrt{1+4\alpha}.$$

**Lemma 9** Assume (A1), (A2), and (A3). Then MLA (2) has local strong error at least of order  $p_2 = 1$ , with maximum step size  $h_2 = \frac{1}{M^2 + 4\alpha}$  and constants

$$C_2 = 7(1+4\alpha) \left( \|y^*\| + \|A(y^*)\|_{HS} + \frac{1}{M} \|g(y^*)\| \right)$$
  
 $D_2 = 5(1+4\alpha).$ 

# 4.3. Proof of Theorem 1: Convergence Rate of MLA

**Proof** [Proof of Theorem 1] Assume (A1), (A2), and (A3) with  $\alpha < m$ . We have verified that MLA satisfies the conditions in the mean-square analysis framework: In Lemma 4 we show MLD is contractive with rate  $\beta = m - \alpha$ . We derive the deviation bound in Lemma 5 with  $C_0 = 4M$ . In Lemmas 8 and 9 we show MLA has local weak error of order  $p_1 = \frac{3}{2}$  and local strong error of order  $p_2 = 1$ , and indeed  $p_2 \le p_1 - \frac{1}{2}$ .

Then by Theorem 3, we can compute the maximum step size:

$$\begin{split} h_{\text{max}} &= \min \left\{ \frac{1}{M^2 + 4\alpha}, \frac{1}{4\beta}, \left( \frac{\sqrt{\beta}}{4\sqrt{2}D_2} \right)^{\frac{1}{p_2 - \frac{1}{2}}}, \left( \frac{\beta}{8\sqrt{2}(D_1 + C_0D_2)} \right)^{\frac{1}{p_2 - \frac{1}{2}}} \right\} \\ &= \min \left\{ \frac{1}{M^2 + 4\alpha}, \frac{1}{4(m - \alpha)}, \frac{m - \alpha}{800(1 + 4\alpha)^2}, \frac{(m - \alpha)^2}{128\left(2M\sqrt{(1 + 4\alpha)} + 20M(1 + 8\alpha)\right)^2} \right\} \\ &= \mathcal{O}\left( \frac{(m - \alpha)^2}{M^2(1 + 4\alpha)^2} \right). \end{split}$$

Recall  $\tilde{\nu} = (\nabla \phi)_{\#} \nu$  is the target distribution of MLD (7). We can compute the constant

$$U = \sqrt{4\mathbb{E}[\|Y_0\|_2^2] + 6\mathbb{E}_{\tilde{\nu}}[\|Y\|_2^2]} = O(\sqrt{d}).$$

Note that 
$$\|A(y^*)\|_{\mathrm{HS}} = \sqrt{\mathrm{Tr}(A(y^*)A(y^*)^\top)} = \sqrt{\mathrm{Tr}(\nabla^2\phi^*(y^*)^{-1})} = \sqrt{\mathrm{Tr}(\nabla^2\phi(x^*))} = O(\sqrt{d}).$$
 Let us define  $V := \|y^*\| + \|A(y^*)\|_{\mathrm{HS}} + \frac{1}{M}\|g(y^*)\| = O(\sqrt{d}).$  Then the resulting constant is 
$$C_{\mathrm{MLA}} = \frac{2}{\beta} \left( C_1 + C_0C_2 + \sqrt{2}U(D_1 + C_0D_2) \right) + \frac{2}{\sqrt{\beta}} \left( C_2 + \sqrt{2}D_2U \right)$$
 
$$= \frac{2}{m-\alpha} \left( 3M\sqrt{(1+4\alpha)}V + 28M(1+4\alpha)V + \sqrt{2}U\left(2M\sqrt{(1+4\alpha)} + 20M(1+4\alpha)\right) \right)$$
 
$$+ \frac{2}{\sqrt{m-\alpha}} \left( 7(1+4\alpha)V + 5\sqrt{2}(1+4\alpha)U \right)$$
 
$$= \mathcal{O}\left( \frac{M(1+4\alpha)\sqrt{d}}{m-\alpha} \right).$$

The conclusion of Theorem 1 follows from Theorem 3.

# 5. Discussion

In this paper, we prove a convergence guarantee for MLA with vanishing bias under modified self-concordance. Our result leaves open many questions, including the following. It would be interesting to consider a more sophisticated discretization of MLD such that the mean-square analysis framework will show improved local errors and smaller bias.

It would be interesting to have a better analysis of MLA under more natural conditions on  $\phi$ , such as self-concordance (rather than modified self-concordance), and under relaxed requirements on f and  $\phi$  (e.g. that allows us to sample from a uniform or Gaussian distribution on a polytope). It would be desirable to have a convergence analysis of MLA in the Wasserstein distance generated by the Hessian metric  $\nabla^2 \phi$  rather than the Euclidean metric, or in other measures such as KL or  $\chi^2$ -divergence.

It would be interesting to understand whether we can discretize the Newton Langevin Dynamics (which is the case when  $\phi=f$  as described in Appendix A.4 and which is affine-invariant in continuous time) and obtain a discrete-time algorithm with a convergence guarantee which is also affine-invariant.

It would also be interesting to understand whether we can derive a more general discrete-time analysis framework that works under a relaxed condition, e.g. without requiring contraction in continuous time, but only exponential convergence in function value (which is known for ULA under the log-Sobolev inequality, see for example Vempala and Wibisono (2019)).

**Acknowledgments.** M.T. is grateful for partial support by NSF DMS–1847802 and ECCS–1936776. S.V. is supported in part by NSF awards CCF–1909756, CCF–2007443 and CCF–2134105.

#### References

Kwangjun Ahn and Sinho Chewi. Efficient constrained sampling via the mirror-Langevin algorithm. *arXiv preprint arXiv:2010.16212*, 2020.

Herm Jan Brascamp and Elliott H Lieb. On extensions of the Brunn-Minkowski and Prékopa-Leindler theorems, including inequalities for log concave functions, and with an application to the diffusion equation. *Journal of Functional Analysis*, 22(4):366–389, 1976. ISSN 0022-1236.

- Sinho Chewi, Thibaut Le Gouic, Chen Lu, Tyler Maunu, Philippe Rigollet, and Austin Stromme. Exponential ergodicity of mirror-Langevin diffusions. In *Advances in Neural Information Processing Systems*, volume 33, pages 19573–19585. Curran Associates, Inc., 2020.
- Arnak Dalalyan. Further and stronger analogy between sampling and optimization: Langevin Monte Carlo and gradient descent. In *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pages 678–689. PMLR, 07–10 Jul 2017. URL http://proceedings.mlr.press/v65/dalalyan17a.html.
- Arnak S Dalalyan and Avetik Karagulyan. User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient. *Stochastic Processes and their Applications*, 2019.
- Alain Durmus and Eric Moulines. Nonasymptotic convergence analysis for the unadjusted Langevin algorithm. *Annals of Applied Probability*, 27(3):1551–1587, 2017.
- Alain Durmus and Eric Moulines. High-dimensional bayesian inference via the unadjusted Langevin algorithm. *Bernoulli*, 25(4A):2854–2882, 2019.
- Max Fathi. *Quelques applications du transport optimal en analyse et en probabilités*. Habilitation à diriger des recherches, Université Paul Sabatier (Toulouse 3), April 2019.
- Ya-Ping Hsieh, Ali Kavis, Paul Rolland, and Volkan Cevher. Mirrored Langevin Dynamics. In Advances in Neural Information Processing Systems 31: NeurIPS 2018, Montréal, Canada, pages 2883–2892, 2018.
- Qijia Jiang. Mirror Langevin Monte Carlo: the case under isoperimetry. In M. Ranzato, A. Beygelzimer, K. Nguyen, P.S. Liang, J.W. Vaughan, and Y. Dauphin, editors, *Advances in Neural Information Processing Systems*, volume 34. Curran Associates, Inc., 2021.
- Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the Fokker–Planck equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, January 1998.
- Ravindran Kannan and Hariharan Narayanan. Random walks on polytopes and an affine interior point method for linear programming. *Mathematics of Operations Research*, 37(1):1–20, 2012. ISSN 0364765X, 15265471.
- Aditi Laddha, Yin Tat Lee, and Santosh Vempala. Strong self-concordance and sampling. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 1212–1222, 2020.
- Ruilin Li, Hongyuan Zha, and Molei Tao. Sqrt(d) dimension dependence of Langevin Monte Carlo. In *International Conference on Learning Representations (ICLR)*, 2022.
- Xuechen Li, Yi Wu, Lester Mackey, and Murat A Erdogdu. Stochastic Runge-Kutta accelerates Langevin Monte Carlo and beyond. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Grigori Noah Milstein and Michael V Tretyakov. *Stochastic numerics for mathematical physics*. Springer Science & Business Media, 2013.

R. Tyrrell Rockafellar. *Convex Analysis*. Princeton Landmarks in Mathematics and Physics. Princeton University Press, 1970. ISBN 978-1-4008-7317-3.

Santosh Vempala and Andre Wibisono. Rapid convergence of the Unadjusted Langevin Algorithm: Isoperimetry suffices. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

Andre Wibisono. Sampling as optimization in the space of measures: The Langevin dynamics as a composite optimization problem. In *Conference On Learning Theory, COLT 2018, Stockholm, Sweden, 6-9 July 2018*, pages 2093–3027, 2018. URL http://proceedings.mlr.press/v75/wibisono18a.html.

Kelvin Shuangjian Zhang, Gabriel Peyré, Jalal Fadili, and Marcelo Pereyra. Wasserstein control of Mirror Langevin Monte Carlo. In *Conference on Learning Theory, COLT 2020, 9-12 July 2020, Virtual Event [Graz, Austria]*, volume 125 of *Proceedings of Machine Learning Research*, pages 3814–3841. PMLR, 2020. URL http://proceedings.mlr.press/v125/zhang20a.html.

# Appendix A. Riemannian and Mirror Langevin Dynamics in Continuous Time

Consider the problem of sampling from  $\nu \propto e^{-f}$  on  $\mathcal{X} \subseteq \mathbb{R}^d$  as described in Section 2.

Suppose we endow  $\mathcal{X}$  with a Riemannian metric g, which we write as a positive definite matrix:  $g(x) \succ 0$  for all  $x \in \mathcal{X}$ . This means at each point  $x \in \mathcal{X}$  we measure local norm using the metric g(x):

$$\langle u, v \rangle_x := u^\top \mathsf{g}(x) v$$

for all u,v in the tangent space. We assume  $x\mapsto \mathsf{g}(x)$  is differentiable. Let  $M(x)=\mathsf{g}(x)^{-1}$  be the inverse matrix, and let  $\sqrt{M(x)}$  be a square-root of M(x). Let  $\nabla\cdot M(x)\in\mathbb{R}^d$  be the divergence of M, which is a vector-valued function whose entries are the divergences of the columns of M. We assume  $\mathsf{g}(x)\to\infty$  (equivalently,  $M(x)\to0$ ) as x approaches the boundary of  $\mathcal X$ .

#### A.1. Review for optimization

Recall in optimization, the **Riemannian gradient flow (RGF)** (or natural gradient flow) for minimizing f using the metric g is the solution  $X_t$  to the differential equation:

$$\dot{X}_t = \frac{d}{dt} X_t = -\mathsf{g}(x)^{-1} \, \nabla f(X_t).$$

Here we use the inverse metric  $M(x) = \mathsf{g}(x)^{-1}$  to turn the  $\ell_2$ -gradient  $\nabla f(x) = (\frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_d})$  into a gradient tangent vector  $\operatorname{grad} f(x) = M(x) \nabla f(x)$  under the Riemannian metric  $\mathsf{g}(x)$ . RGF has nice properties when the objective function f satisfies some properties. For example, if f is  $\operatorname{geodesically strongly convex}$  (which means f is strongly convex along geodesics generated by the Riemannian metric  $\mathsf{g}$ ), then RGF is exponentially contracting. Moreover, if f is  $\operatorname{gradient dominated}$  with respect to  $\mathsf{g}$ , then the function value  $f(X_t)$  converges exponentially fast along RGF.

Consider when the metric g(x) is given by the Hessian of a convex Legendre function  $\phi$ :  $g(x) = \nabla^2 \phi(x) > 0$ . Then the RGF becomes:

$$\dot{X}_t = -\nabla^2 \phi(X_t)^{-1} \, \nabla f(X_t).$$

In terms of the dual variable  $Y_t = \nabla \phi(X_t)$ , this becomes the **mirror flow**:

$$\dot{Y}_t = -\nabla f(X_t) = -\nabla f(\nabla \phi^*(Y_t)).$$

Recall by the mirror map  $\nabla \phi$ , the metric  $\nabla^2 \phi$  on  $\mathcal X$  becomes the Hessian metric  $\nabla^2 \phi^*$  on  $\mathcal Y = \nabla \phi(X) = \mathbb R^d$ . The mirror flow is also the Riemannian gradient flow for minimizing the push-forward function  $\tilde f(y) = f(\nabla \phi^*(y))$  under the Hessian metric  $\nabla^2 \phi^*(y)$  (because  $\operatorname{grad} \tilde f(y) = \nabla^2 \phi^*(y)^{-1} \nabla \tilde f(y) = \nabla f(\nabla \phi^*(y))$ ). Discretizing the mirror flow gives the *mirror descent* algorithm in optimization.

# A.2. Riemannian Langevin Dynamics

The **Riemannian Langevin Dynamics (RLD)** for sampling from  $\nu \propto e^{-f}$  on  $\mathcal{X}$  using the metric g(x) is the solution  $X_t$  to the stochastic differential equation:

$$dX_t = (\nabla \cdot M(X_t) - M(X_t) \nabla f(X_t)) dt + \sqrt{2} \sqrt{M(X_t)} dW_t$$
(9)

where  $M(x) = g(x)^{-1}$ . Here  $W_t$  is the standard Brownian motion in  $\mathbb{R}^d$ . Since  $M(x) \to 0$  as  $x \to \partial \mathcal{X}$ , the process does not leave  $\mathcal{X}$ : If  $X_0 \in \mathcal{X}$ , then  $X_t \in \mathcal{X}$  for all t > 0.

The additional drift term  $\nabla \cdot M(X_t)$  accounts for the covariance  $M(X_t)$  in the Brownian motion. The stationary distribution for RLD is  $\nu(x) \propto e^{-f(x)}$  (the density is with respect to the Lebesgue measure dx on  $\mathbb{R}^d$ ). This can be seen, for example, from the following Fokker-Planck equation. If  $X_t \in \mathcal{X}$  follows RLD (9), then its density  $\rho_t \colon \mathcal{X} \to \mathbb{R}$  evolves following the partial differential equation (PDE):

$$\frac{\partial \rho_t}{\partial t} = \nabla \cdot \left( \rho_t M \nabla \log \frac{\rho_t}{\nu} \right). \tag{10}$$

Clearly if  $\rho_t = \nu$  then the dynamics is stationary. Furthermore, the PDE above can be interpreted as the gradient flow for minimizing relative entropy with respect to the Wasserstein metric on the metric space  $(\mathcal{X}, g)$ .

From the Fokker-Planck equation (10), we can derive how fast the dynamics RLD approaches the target distribution  $\nu$  in various measures.

For example, recall the  $\chi^2$ -divergence of a probability distribution  $\rho$  with respect to  $\nu$  is

$$\chi_{\nu}^{2}(\rho) = \operatorname{Var}_{\nu}\left(\frac{\rho}{\nu}\right) = \int_{\mathcal{X}} \nu(x) \left(\frac{\rho(x)}{\nu(x)} - 1\right)^{2} dx = \int_{\mathcal{X}} \frac{\rho(x)^{2}}{\nu(x)} dx - 1.$$

Then a standard calculation reveals that the  $\chi^2$ -divergence is decreasing along RLD (10):

$$\frac{d}{dt}\chi_{\nu}^{2}(\rho_{t}) = -2G_{\nu}(\rho_{t})$$

where

$$G_{\nu}(\rho) = \mathbb{E}_{\nu} \left[ \left\| \nabla \left( \frac{\rho}{\nu} \right) \right\|_{M}^{2} \right] = \int_{\mathcal{X}} \nu(x) \left\langle \nabla \left( \frac{\rho(x)}{\nu(x)} \right), M(x) \nabla \left( \frac{\rho(x)}{\nu(x)} \right) \right\rangle dx.$$

Therefore, if  $\nu$  satisfies a *Poincaré inequality* with respect to g, which means for any differentiable function  $h \colon \mathcal{X} \to \mathbb{R}$  we have

$$\operatorname{Var}_{\nu}(h) \leq C_{\operatorname{P}} \mathbb{E}_{\nu}[\|\nabla h\|_{M}^{2}],$$

then we can conclude RLD converges exponentially fast in  $\chi^2$ -divergence:  $\chi^2_{\nu}(\rho_t) \leq e^{-\frac{2t}{C_{\rm P}}} \chi^2_{\nu}(\rho_0)$ . Similarly, recall the *relative entropy* (or *KL divergence*) of  $\rho$  with respect to  $\nu$  is

$$H_{\nu}(\rho) = \mathbb{E}_{\nu} \left[ \frac{\rho}{\nu} \log \frac{\rho}{\nu} \right] = \int_{\mathcal{X}} \rho(x) \log \frac{\rho(x)}{\nu(x)} dx.$$

Then along RLD (10), KL divergence is decreasing:

$$\frac{d}{dt}H_{\nu}(\rho_t) = -J_{\nu}(\rho_t)$$

where  $J_{\nu}(\rho)$  is the relative Fisher information of  $\rho$  with respect to  $\nu$  under the metric g:

$$G_{\nu}(\rho) = \mathbb{E}_{\rho} \left[ \left\| \nabla \log \frac{\rho}{\nu} \right\|_{M}^{2} \right].$$

Therefore, if  $\nu$  satisfies a log-Sobolev inequality with respect to g, which means for any  $\rho$  we have

$$H_{\nu}(\rho) \leq C_{\text{LSI}} J_{\nu}(\rho),$$

then we can conclude RLD converges exponentially fast in KL divergence:  $H_{\nu}(\rho_t) \leq e^{-\frac{t}{C_{\rm LSI}}} H_{\nu}(\rho_0)$ .

# A.3. Mirror Langevin Dynamics

Suppose now the metric g(x) is given by the Hessian of a convex Legendre function  $\phi$ :  $g(x) = \nabla^2 \phi(x) > 0$ . The Riemannian Langevin dynamics (9) becomes the following SDE, which is also studied by Zhang et al. (2020) and Chewi et al. (2020):

$$dX_t = \left(\nabla \cdot (\nabla^2 \phi(X_t)^{-1}) - \nabla^2 \phi(X_t)^{-1} \nabla f(X_t)\right) dt + \sqrt{2\nabla^2 \phi(X_t)^{-1}} dW_t.$$
 (11)

If  $\nu$  satisfies log-Sobolev or Poincaré inequality (which is called *mirror Poincaré inequality* in Chewi et al. (2020)), then we can conclude exponential convergence rate in KL or  $\chi^2$  divergence along (11). The SDE (11) requires  $\nabla \cdot (\nabla^2 \phi(x)^{-1})$ , which may be complicated. Consider the dual variable  $Y_t = \nabla \phi(X_t)$ . By Itô's lemma,  $Y_t$  evolves following the **Mirror Langevin Dynamics**:

$$dY_t = -\nabla f(\nabla \phi^*(Y_t)) dt + \sqrt{2\nabla^2 \phi^*(Y_t)^{-1}} dW_t.$$

For an explicit calculation, see for example (Jiang, 2021, Appendix A). In particular, the drift term simplifies and there is no third derivative involved. The mirror Langevin dynamics is also the Riemannian Langevin dynamics (9) for sampling from the pushforward distribution  $\tilde{\nu} = (\nabla \phi)_{\#} \nu$  using the Hessian metric  $\nabla^2 \phi^*$ . Furthermore, the  $\chi^2$ -divergence and KL divergence are invariant under the mirror map. Therefore,  $\nu$  satisfies LSI or Poincaré inequality with respect to  $\nabla^2 \phi$  if and only if  $\tilde{\nu}$  also satisfies LSI or Poincaré inequality with respect to  $\nabla^2 \phi^*$ . Therefore, we get the same convergence guarantee in both primal and dual spaces.

# A.4. Newton Langevin Dynamics

A particularly nice choice of  $\phi$  is when  $\phi = f$ . This gives the **Newton Langevin Dynamics**, which in the primal space takes the form:

$$dX_t = \left(\nabla \cdot (\nabla^2 f(X_t)^{-1}) - \nabla^2 f(X_t)^{-1} \nabla f(X_t)\right) dt + \sqrt{2} \sqrt{\nabla^2 f(X_t)^{-1}} dW_t. \tag{12}$$

A remarkable property of NLD, as pointed out by Chewi et al. (2020), is that the Poincaré inequality of  $\nu \propto e^{-f}$  with respect to its Hessian metric  $\nabla^2 f$  is always true with a uniform constant  $C_P = 1$  for any strictly log-concave distribution  $\nu$ , by the virtue of the Brascamp-Lieb inequality (Brascamp and Lieb, 1976). This gives a uniform exponential convergence rate along NLD in  $\chi^2$ -divergence as well as the Wasserstein distance with respect to the metric  $\nabla^2 f$ ; see detailed exposition and additional consequences in Chewi et al. (2020).

In the dual space, Newton Langevin Dynamics has a simple drift:

$$dY_t = -Y_t dt + \sqrt{2\nabla^2 f^*(Y_t)^{-1}} dW_t$$
(13)

since  $\nabla f(\nabla f^*(y)) = y$ . The target distribution of NLD in the dual space is the pushforward distribution  $\tilde{\nu} = (\nabla f)_{\#} \nu$  where  $\nu \propto e^{-f}$ . The SDE (13) for sampling from  $\tilde{\nu}$  was also pointed out by Fathi (2019) from the study of Stein's kernel.

# Appendix B. Proofs of Lemmas

#### B.1. Proof of Lemma 4: Contraction of MLD

**Proof** [Proof of Lemma 4] Assume (A1) and (A2). We will show MLD (7) is contractive if  $\alpha < \frac{m}{2}$ . Suppose we have two solutions  $Y'_t, Y_t$  of MLD (7) with the same Brownian motion:

$$dY'_t = -g(Y'_t)dt + \sqrt{2}A(Y'_t)dW_t$$
  
$$dY_t = -g(Y_t)dt + \sqrt{2}A(Y_t)dW_t.$$

Then the difference satisfies the SDE

$$d(Y_t' - Y_t) = -(g(Y_t') - g(Y_t))dt + \sqrt{2}(A(Y_t') - A(Y_t))dW_t.$$
(14)

Recall in general that if  $V_t \in \mathbb{R}^d$  follows a general SDE  $dV_t = b(V_t)dt + G(V_t)dW_t$ , then

$$\frac{d}{dt}\mathbb{E}[\|V_t\|^2] = \mathbb{E}[2\langle b(V_t), V_t \rangle + \|G(V_t)\|_{HS}^2].$$

Then for the SDE (14) of the difference  $V_t = Y'_t - Y_t$ , and by applying assumptions (A1) and (A3), we have

$$\frac{d}{dt}\mathbb{E}[\|Y'_t - Y_t\|^2] = -2\mathbb{E}[\langle g(Y'_t) - g(Y_t), Y'_t - Y_t \rangle] + 2\mathbb{E}[\|A(Y'_t) - A(Y_t)\|_{HS}^2] 
\leq -2m\mathbb{E}[\|Y'_t - Y_t\|_2^2] + 2\alpha\mathbb{E}[\|Y'_t - Y_t\|_2^2] 
= -2(m - \alpha)\mathbb{E}[\|Y'_t - Y_t\|_2^2].$$

We see that we have an exponential contraction if  $\alpha < m$ :

$$\mathbb{E}[\|Y_t' - Y_t\|^2] \le \exp(-2(m - \alpha)t) \,\mathbb{E}[\|Y_0' - Y_0\|^2] \quad \forall \, t \ge 0.$$
 (15)

This shows that MLD (7) is contractive with rate  $\beta = m - \alpha$ .

#### B.2. Proof of Lemma 5: Deviation bound of MLD

**Proof** [Proof of Lemma 5] Assume (A1), (A2), and (A3) with  $\alpha < m$ .

Suppose we have two solutions  $Y'_t, Y_t$  of MLD (7) with the same Brownian motion:

$$dY'_t = -g(Y'_t)dt + \sqrt{2}A(Y'_t)dW_t$$
  
$$dY_t = -g(Y_t)dt + \sqrt{2}A(Y_t)dW_t.$$

Consider the shifted variables  $\tilde{Y}'_t = Y'_t - Y'_0$  and  $\tilde{Y}_t = Y_t - Y_0$ , which satisfy:

$$d\tilde{Y}'_{t} = -g(\tilde{Y}'_{t} + Y'_{0})dt + \sqrt{2}A(\tilde{Y}'_{t} + Y'_{0})dW_{t}$$
  
$$d\tilde{Y}_{t} = -g(\tilde{Y}_{t} + Y_{0})dt + \sqrt{2}A(\tilde{Y}_{t} + Y_{0})dW_{t}.$$

Then the difference  $\tilde{Y}_t' - \tilde{Y}_t = (Y_t' - Y_0') - (Y_t - Y_0)$  satisfies:

$$d(\tilde{Y}_t' - \tilde{Y}_t) = -(g(\tilde{Y}_t' + Y_0') - g(\tilde{Y}_t + Y_0))dt + \sqrt{2}(A(\tilde{Y}_t' + Y_0') - A(\tilde{Y}_t + Y_0))dW_t.$$

By Lemma 4, we have the contraction result (15), which implies  $\mathbb{E}[\|Y_t' - Y_t\|_2^2] \leq \mathbb{E}[\|Y_0' - Y_0\|_2^2]$  for all  $t \geq 0$ . Then by applying (A1) and (A2) and using  $\alpha < m \leq M$ , we get

$$\begin{split} &\frac{d}{dt}\mathbb{E}[\|(Y_t'-Y_0')-(Y_t-Y_0)\|_2^2] \\ &=\frac{d}{dt}\mathbb{E}[\|\tilde{Y}_t'-\tilde{Y}_t\|_2^2] \\ &=-2\mathbb{E}[\langle g(\tilde{Y}_t'+Y_0')-g(\tilde{Y}_t+Y_0),\tilde{Y}_t'-\tilde{Y}_t\rangle]+2\mathbb{E}[\|A(\tilde{Y}_t'+Y_0')-A(\tilde{Y}_t+Y_0)\|_{\mathrm{HS}}^2] \\ &=-2\mathbb{E}[\langle g(Y_t')-g(Y_t),Y_t'-Y_t-(Y_0'-Y_0)\rangle]+2\mathbb{E}[\|A(Y_t')-A(Y_t)\|_{\mathrm{HS}}^2] \\ &\leq 2\mathbb{E}[\langle g(Y_t')-g(Y_t),Y_0'-Y_0\rangle]+2\alpha\mathbb{E}[\|Y_t'-Y_t\|_2^2] \\ &\leq 2\mathbb{E}[\|g(Y_t')-g(Y_t)\|_2^2]^{1/2}\mathbb{E}[\|Y_0'-Y_0\|_2^2]^{\frac{1}{2}}+2\alpha\mathbb{E}[\|Y_0'-Y_0\|_2^2] \\ &\leq 2M\mathbb{E}[\|Y_t'-Y_t\|_2^2]^{\frac{1}{2}}\mathbb{E}[\|Y_0'-Y_0\|_2^2]^{\frac{1}{2}}+2\alpha\mathbb{E}[\|Y_0'-Y_0\|_2^2] \\ &\leq (2M+2\alpha)\,\mathbb{E}[\|Y_0'-Y_0\|_2^2] \\ &\leq 4M\,\mathbb{E}[\|Y_0'-Y_0\|_2^2]. \end{split}$$

Integrating, we conclude that for all  $t \ge 0$ ,

$$\mathbb{E}[\|(Y_t' - Y_0') - (Y_t - Y_0)\|_2^2] \le 4M \, \mathbb{E}[\|Y_0' - Y_0\|_2^2] \, t.$$

# B.3. Proof of Lemma 6: Growth bound of MLD

**Proof** [Proof of Lemma 6] Assume (A1), (A2), and (A3). Consider the solution  $Y_t$  of MLD (7) starting from  $Y_0$ . The centered variable  $\tilde{Y}_t = Y_t - Y_0$  follows the SDE

$$d\tilde{Y}_t = -g(\tilde{Y}_t + Y_0)dt + \sqrt{2}A(\tilde{Y}_t + Y_0)dW_t.$$

Then

$$\frac{d}{dt}\mathbb{E}[\|Y_t - Y_0\|_2^2] = \frac{d}{dt}\mathbb{E}[\|\tilde{Y}_t\|_2^2] = -2\mathbb{E}[\langle g(\tilde{Y}_t + Y_0), \tilde{Y}_t \rangle] + 2\mathbb{E}[\|A(\tilde{Y}_t + Y_0)\|_{\mathrm{HS}}^2]$$

$$= \underbrace{-2\mathbb{E}[\langle g(Y_t), Y_t - Y_0 \rangle]}_{=I} + \underbrace{2\mathbb{E}[\|A(Y_t)\|_{\mathrm{HS}}^2]}_{=II}.$$

Let us bound the two terms above. Let  $x^* = \arg\min_{x \in \mathcal{X}} f(x)$  and  $y^* = \nabla \phi(x^*)$ .

First term: By (A2) and (A3),

$$\begin{split} I &= -2\mathbb{E}[\langle g(Y_t), Y_t - Y_0 \rangle] \\ &= -2\mathbb{E}[\langle g(Y_t) - g(Y_0), Y_t - Y_0 \rangle] - 2\mathbb{E}[\langle g(Y_0), Y_t - Y_0 \rangle] \\ &\leq -2\mathbb{E}[\langle g(Y_0), Y_t - Y_0 \rangle] \\ &\leq 2\mathbb{E}[\|g(Y_0)\| \cdot \|Y_t - Y_0\|] \\ &\leq \frac{1}{M^2}\mathbb{E}[\|g(Y_0)\|^2] + M^2\mathbb{E}[\|Y_t - Y_0\|^2] \\ &\leq 2\mathbb{E}[\|Y_0 - y^*\|^2] + M^2\mathbb{E}[\|Y_t - Y_0\|^2] + \frac{2}{M^2}\|g(y^*)\|^2. \end{split}$$

In the last step we have used  $||g(y)||_2^2 \le 2||g(y) - g(y^*)||^2 + 2||g(y^*)||^2 \le 2M^2||y - y^*||^2 + 2||g(y^*)||^2$ .

**Second term:** By triangle inequality and (A1),

$$||A(Y_t)||_{\mathrm{HS}}^2 \le 2||A(Y_t) - A(Y_0)||_{\mathrm{HS}}^2 + 2||A(Y_0)||_{\mathrm{HS}}^2$$

$$\le 2||A(Y_t) - A(Y_0)||_{\mathrm{HS}}^2 + 4||A(Y_0) - A(y^*)||_{\mathrm{HS}}^2 + 4||A(y^*)||_{\mathrm{HS}}^2$$

$$\le 2\alpha ||Y_t - Y_0||_2^2 + 4\alpha ||Y_0 - y^*||^2 + 4||A(y^*)||_{\mathrm{HS}}^2.$$

Therefore,

$$II = 2\mathbb{E}[\|A(Y_t)\|_{HS}^2]$$
  

$$\leq 4\alpha \mathbb{E}[\|Y_t - Y_0\|_2^2] + 8\alpha \mathbb{E}[\|Y_0 - y^*\|^2] + 8\|A(y^*)\|_{HS}^2.$$

Combining the two terms above, we get that along MLD (7):

$$\frac{d}{dt}\mathbb{E}[\|Y_t - Y_0\|_2^2] \le (M^2 + 4\alpha)\mathbb{E}[\|Y_t - Y_0\|_2^2] + D \tag{16}$$

where

$$D = (2 + 8\alpha)\mathbb{E}[\|Y_0 - y^*\|^2] + 8\|A(y^*)\|_{HS}^2 + \frac{2}{M^2}\|g(y^*)\|^2$$

$$\leq 4(1 + 4\alpha)\mathbb{E}[\|Y_0\|^2] + 4(1 + 4\alpha)\|y^*\|^2 + 8\|A(y^*)\|_{HS}^2 + \frac{2}{M^2}\|g(y^*)\|^2.$$

Recall in general if  $V_t \geq 0$  satisfies  $\frac{d}{dt}V_t \leq CV_t + D$  for some C, D > 0, then  $V_t \leq e^{Ct}V_0 + \frac{D}{C}(e^{Ct}-1)$ . Furthermore, if  $V_0 = 0$  and  $0 < t \leq \frac{1}{C}$ , then  $V_t \leq \frac{D}{C}2Ct = 2Dt$ . Applying this to  $V_t = \mathbb{E}[\|Y_t - Y_0\|^2]$  which satisfies (16) and  $V_0 = 0$ , we conclude that if  $0 < t \leq \frac{1}{M^2 + 4\alpha}$ , then

$$\mathbb{E}[\|Y_t - Y_0\|_2^2] \le \gamma t$$

where 
$$\gamma = 2D \le 8(1+4\alpha)\mathbb{E}[\|Y_0\|^2] + 8(1+4\alpha)\|y^*\|^2 + 16\|A(y^*)\|_{\mathrm{HS}}^2 + \frac{4}{M^2}\|g(y^*)\|^2$$
.

#### B.4. Proof of Lemma 8: Local weak error of MLA

**Proof** [Proof of Lemma 8] Assume (A1), (A2), and (A3). Starting from  $Y_0 \in \mathbb{R}^d$ , let  $Y_t$  be the solution to the MLD (7), and let  $Y_t'$  be the solution to the modified SDE with constant drift, driven by the same Brownian motion:

$$dY_t = -g(Y_t)dt + \sqrt{2}A(Y_t)dW_t$$
  
$$dY_t' = -g(Y_0)dt + \sqrt{2}A(Y_0)dW_t.$$

The value  $Y'_h$  at time t = h is the output  $\bar{Y}_1$  of MLA (2) from  $Y_0$ . We wish to bound  $\|\mathbb{E}[Y_h - Y'_h]\|$ . Since  $Y_t, Y'_t$  are coupled using the same Brownian motion, the difference  $Y_t - Y'_t$  satisfies

$$d(Y_t - Y_t') = -(g(Y_t) - g(Y_0))dt + \sqrt{2}(A(Y_t) - A(Y_0))dW_t.$$

Integrating, and since  $Y_0 = Y_0'$ , this means

$$Y_h - Y_h' = -\int_0^h (g(Y_t) - g(Y_0))dt + \sqrt{2} \int_0^h (A(Y_t) - A(Y_0))dW_t.$$

Taking expectation gives

$$\mathbb{E}[Y_h - Y_h'] = -\int_0^h \mathbb{E}[g(Y_t) - g(Y_0)]dt. \tag{17}$$

By (A2) and Lemma 6, for  $0 < t \le \frac{1}{M^2 + 4\alpha}$  we have

$$\mathbb{E}[\|g(Y_t) - g(Y_0)\|] \le M \mathbb{E}[\|Y_t - Y_0\|] \le M \sqrt{\mathbb{E}[\|Y_t - Y_0\|^2]} \le M \sqrt{\gamma t}.$$

Therefore, by triangle inequality on (17), for  $0 < h \le \frac{1}{M^2 + 4\alpha}$  we have

$$\begin{split} \|\mathbb{E}[Y_h - Y_h']\| &\leq \int_0^h \mathbb{E}[\|g(Y_t) - g(Y_0)\|] dt \\ &\leq M \sqrt{\gamma} \int_0^h \sqrt{t} \, dt \\ &= \frac{2}{3} M \sqrt{\gamma} \, h^{\frac{3}{2}} \\ &= \frac{2}{3} M \left( 8(1 + 4\alpha) \mathbb{E}[\|Y_0\|^2] + 8(1 + 4\alpha) \|y^*\|^2 + 16 \|A(y^*)\|_{\mathrm{HS}}^2 + \frac{4}{M^2} \|g(y^*)\|^2 \right)^{\frac{1}{2}} h^{\frac{3}{2}} \\ &\leq \frac{2}{3} M \left( \sqrt{8(1 + 4\alpha)} \sqrt{\mathbb{E}[\|Y_0\|^2]} + \sqrt{8(1 + 4\alpha)} \|y^*\| + 4 \|A(y^*)\|_{\mathrm{HS}} + \frac{2}{M} \|g(y^*)\| \right) h^{\frac{3}{2}} \\ &= \left( C_1 + D_1 \sqrt{\mathbb{E}[\|Y_0\|^2]} \right) h^{3/2}. \end{split}$$

This shows the local weak error order is at least  $p_1 = \frac{3}{2}$ , with maximum step size  $h_1 = \frac{1}{M^2 + 8\alpha}$  and constants

$$C_{1} = \frac{2}{3}M\left(\sqrt{8(1+4\alpha)}\|y^{*}\| + 4\|A(y^{*})\|_{HS} + \frac{2}{M}\|g(y^{*})\|\right)$$

$$\leq 3M\sqrt{(1+4\alpha)}\left(\|y^{*}\| + \|A(y^{*})\|_{HS} + \frac{1}{M}\|g(y^{*})\|\right)$$

$$D_{1} = \frac{2}{3}M\sqrt{8(1+4\alpha)}$$

$$\leq 2M\sqrt{1+4\alpha}.$$

# B.5. Proof of Lemma 9: Local strong error of MLA

**Proof** [Proof of Lemma 9]

Assume (A1), (A2), and (A3). As in the proof of Lemma 8, consider two dynamics  $Y'_t, Y_t$  starting from  $Y'_0 = Y_0$  following the SDEs coupled with the same Brownian motion:

$$dY_t = -g(Y_t)dt + \sqrt{2}A(Y_t)dW_t$$
  
$$dY_t' = -g(Y_0)dt + \sqrt{2}A(Y_0)dW_t.$$

We wish to bound  $\mathbb{E}[||Y_h - Y_h'||^2]$ . The difference  $Y_t - Y_t'$  satisfies

$$d(Y_t - Y_t') = -(g(Y_t) - g(Y_0))dt + \sqrt{2}(A(Y_t) - A(Y_0))dW_t.$$

By (A1), (A2), and Lemma 6, for  $0 < t \le \frac{1}{M^2 + 4\alpha}$  we have

$$\begin{split} \frac{d}{dt} \mathbb{E}[\|Y_t - Y_t'\|_2^2] &= -2\mathbb{E}[\langle g(Y_t) - g(Y_0), Y_t - Y_t' \rangle] + 2\mathbb{E}[\|A(Y_t) - A(Y_0)\|_{\mathrm{HS}}^2] \\ &\leq 2\mathbb{E}[\|g(Y_t) - g(Y_0)\|^2]^{\frac{1}{2}} \, \mathbb{E}[\|Y_t - Y_t'\|^2]^{\frac{1}{2}} + 2\alpha \mathbb{E}[\|Y_t - Y_0\|^2] \\ &\leq 2M \mathbb{E}[\|Y_t - Y_0\|^2]^{\frac{1}{2}} \, \mathbb{E}[\|Y_t - Y_t'\|^2]^{\frac{1}{2}} + 2\alpha \mathbb{E}[\|Y_t - Y_0\|^2] \\ &\leq M^2 \mathbb{E}[\|Y_t - Y_t'\|^2] + (1 + 2\alpha)\mathbb{E}[\|Y_t - Y_0\|^2] \\ &\leq M^2 \mathbb{E}[\|Y_t - Y_t'\|^2] + (1 + 2\alpha)\gamma t. \end{split}$$

Equivalently,  $\frac{d}{dt}(e^{-M^2t}\mathbb{E}[||Y_t - Y_t'||_2^2]) \le e^{-M^2t}(1 + 2\alpha)\gamma t \le (1 + 2\alpha)\gamma t$ , so

$$\mathbb{E}[\|Y_t - Y_t'\|_2^2] \le e^{M^2 t} \frac{(1+2\alpha)}{2} \gamma t^2.$$

Furthermore, since  $t \leq \frac{1}{M^2 + 4\alpha} \leq \frac{1}{M^2}$ , we have  $e^{M^2 t} \leq e < 3$ , so

$$\mathbb{E}[\|Y_t - Y_t'\|_2^2] \le \frac{3}{2}(1 + 2\alpha)\gamma t^2$$

$$= 3(1 + 2\alpha)\left(8(1 + 4\alpha)\mathbb{E}[\|Y_0\|^2] + 8(1 + 4\alpha)\|y^*\|^2 + 16\|A(y^*)\|_{HS}^2 + \frac{4}{M^2}\|g(y^*)\|^2\right)t^2$$

$$= (C_2^2 + D_2^2 \mathbb{E}[\|Y_0\|^2])t^2.$$

This shows the local strong error order is at least  $p_2 = 1$  with maximum step size  $h_2 = \frac{1}{M^2 + 4\alpha}$  and constants

$$C_{2} = \left(24(1+2\alpha)(1+4\alpha)\|y^{*}\|^{2} + 48(1+2\alpha)\|A(y^{*})\|_{HS}^{2} + \frac{12(1+2\alpha)}{M^{2}}\|g(y^{*})\|^{2}\right)^{\frac{1}{2}}$$

$$\leq 5(1+4\alpha)\|y^{*}\| + 7\sqrt{1+2\alpha}\|A(y^{*})\|_{HS} + \frac{4\sqrt{1+2\alpha}}{M}\|g(y^{*})\|$$

$$\leq 7(1+4\alpha)\left(\|y^{*}\| + \|A(y^{*})\|_{HS} + \frac{1}{M}\|g(y^{*})\|\right)$$

$$D_{2} = \sqrt{24(1+2\alpha)(1+4\alpha)}$$

$$\leq 5(1+4\alpha).$$

# Appendix C. An Analogy: Geometric Brownian Motion

We wondered if our requirement on the modified self-concordance parameter  $\alpha$  being upper-bounded is an artifact of our proof technique. Thus we did some simple calculations on **Geometric Brownian Motion (GBM)** which is an SDE with multiplicative noise and yet admitting close-form solution. It is not an exact example of MLD but only an analogy; nevertheless, GBM does need  $\alpha$  to be bounded in order to converge.

More precisely, consider GBM on  $\mathbb{R}_+ = (0, \infty)$  which follows the stochastic differential equation:

$$dY_t = -Y_t dt + \sqrt{2\alpha} Y_t dW_t \tag{18}$$

where  $dW_t$  is the standard Brownian motion on  $\mathbb{R}$ . This has exact solution

$$Y_t = Y_0 \exp\left(-(1+\alpha)t + \sqrt{2\alpha} W_t\right).$$

By a standard calculation, we see there is a threshold  $\alpha < 1$  for the convergence of  $Y_t$  as  $t \to \infty$ . Recall since  $W_t \sim \mathcal{N}(0,t)$ ,  $\mathbb{E}[\exp(\sigma W_t)] = e^{\sigma^2 t/2}$  for all  $\sigma > 0$ . Then

$$\mathbb{E}[Y_t^2] = \mathbb{E}[Y_0^2] \, e^{-2(1+\alpha)t} \, \mathbb{E}[\exp(2\sqrt{2\alpha}W_t)] = \mathbb{E}[Y_0^2] \, e^{-2(1-\alpha)t}.$$

Therefore,

$$\lim_{t \to \infty} \mathbb{E}[Y_t^2] = \begin{cases} 0 & \text{if } \alpha < 1 \\ \mathbb{E}[Y_0^2] & \text{if } \alpha = 1 \\ \infty & \text{if } \alpha > 1. \end{cases}$$

Now consider a synchronous coupling  $Y_t, \tilde{Y}_t$  following GBM (18) with the same Brownian motion:

$$Y_t = Y_0 \exp\left(-(1+\alpha)t + \sqrt{2\alpha} W_t\right)$$
  
$$\tilde{Y}_t = \tilde{Y}_0 \exp\left(-(1+\alpha)t + \sqrt{2\alpha} W_t\right).$$

Then

$$\mathbb{E}[(Y_t - \tilde{Y}_t)^2] = \mathbb{E}[(Y_0 - \tilde{Y}_0)^2] e^{-2(1+\alpha)t} \mathbb{E}[\exp(2\sqrt{2\alpha}W_t)] = \mathbb{E}[(Y_0 - \tilde{Y}_0)^2] e^{-2(1-\alpha)t}.$$

Thus, we see that GBM is a contraction if and only if  $\alpha < 1$ . In particular, we also have

$$\lim_{t \to \infty} \mathbb{E}[(Y_t - \tilde{Y}_t)^2] = \begin{cases} 0 & \text{if } \alpha < 1\\ \mathbb{E}[(Y_0 - \tilde{Y}_0)^2] & \text{if } \alpha = 1\\ \infty & \text{if } \alpha > 1. \end{cases}$$

GBM (18) is an instance of MLD (7) (and in fact NLD (12)) with  $\phi = f$  where

$$\frac{1}{\sqrt{(\phi^*)''(y)}} = \sqrt{\alpha}y\tag{19}$$

which is  $\sqrt{\alpha}$ -Lipschitz, so it satisfies modified self-concordance (A1) with parameter  $\alpha$ . Since  $\phi=f$ , it satisfes relative smoothness (A2) and relative strong convexity (A3) with M=m=1. Note our assumption in Theorem 1 is  $\alpha < m=1$ , which c is tight for GBM to contract, as well as to determine if there is a  $t\to\infty$  limit.

# Appendix D. Example: Log-Barrier on a Polytope

Let  $\mathcal{X}$  be the polytope (not necessarily bounded)

$$\mathcal{X} = \{ x \in \mathbb{R}^d \colon a_i^\top x \ge b_i \ \forall i = 1, \dots, m \}$$

for some  $a_1, \ldots, a_m \in \mathbb{R}^d$  and  $b_1, \ldots, b_m \in \mathbb{R}$ . Consider the log-barrier function defined in the interior of  $\mathcal{X}$ :

$$\phi(x) = -\sum_{i=1}^{m} \log(a_i^{\top} x - b_i).$$

Recall that  $\phi$  satisfies the classical self-concordance condition with a constant parameter 2. Let  $\alpha$  be the modified self-concordance parameter of  $\phi$ . For some polytopes, such as the positive orthant,  $\alpha$  is also a constant (because the Hessian is diagonal and the dimensions are independent). For general polytopes, however,  $\alpha$  can be arbitrarily large. Here we show  $\alpha$  can be as large as the square inverse of the smallest singular value of the constraint matrix; we also construct an explicit example in two dimension.

Without loss of generality we may assume  $\|a_i\|=1$  for  $i=1,\ldots,m$ . Let  $A=(a_1,\cdots,a_m)\in\mathbb{R}^{d\times m}$  be the constraint matrix, so the polytope is described by  $A^\top x\geq b$ . Let the singular values of A be  $\sigma_1\geq\cdots\geq\sigma_d\geq0$  (assuming  $d\leq m$ ). Then  $\sum_{i=1}^d\sigma_i^2=\operatorname{Tr}(AA^\top)=\operatorname{Tr}(A^\top A)=\sum_{i=1}^m\|a_i\|^2=m$ ; but  $\sigma_d=\min_i\sigma_i$  can be small or 0. For  $x\in\mathcal{X}$ , let  $S_x\in\mathbb{R}^{m\times m}$  be the diagonal matrix with entries  $a_i^\top x-b_i$ .

The gradient of  $\phi$  is

$$\nabla \phi(x) = -\sum_{i=1}^{m} \frac{a_i}{a_i^{\top} x - b_i}.$$

Then for  $x, x' \in \mathcal{X}$ , we have

$$\nabla \phi(x') - \nabla \phi(x) = \sum_{i=1}^{m} \left( \frac{1}{a_i^{\top} x - b_i} - \frac{1}{a_i^{\top} x' - b_i} \right) a_i$$

$$= \sum_{i=1}^{m} \frac{a_i^{\top} (x' - x)}{(a_i^{\top} x - b_i)(a_i^{\top} x' - b_i)} a_i$$

$$= \sum_{i=1}^{m} \frac{a_i a_i^{\top}}{(a_i^{\top} x - b_i)(a_i^{\top} x' - b_i)} (x' - x)$$

$$= A S_x^{-1} S_{x'}^{-1} A^{\top} (x' - x).$$

Therefore,

$$\|\nabla\phi(x') - \nabla\phi(x)\|^2 = \|AS_x^{-1}S_{x'}^{-1}A^\top(x'-x)\|^2$$

$$= (x'-x)^\top AS_{x'}^{-1}S_x^{-1}A^\top AS_x^{-1}S_{x'}^{-1}A^\top(x'-x)$$

$$= v(x,x')^\top A^\top Av(x,x')$$

where

$$v(x, x') = S_x^{-1} S_{x'}^{-1} A^{\top} (x' - x) \in \mathbb{R}^m.$$

The Hessian is

$$\nabla^2 \phi(x) = \sum_{i=1}^m \frac{a_i a_i^{\top}}{(a_i^{\top} x - b_i)^2} = A S_x^{-2} A^{\top}.$$

As a square-root, we can choose:

$$\sqrt{\nabla^2 \phi(x)} = AS_x^{-1} = \left(\frac{a_1}{a_1^\top x - b_1} \quad \cdots \quad \frac{a_m}{a_m^\top x - b_m}\right)$$

since indeed  $\sqrt{\nabla^2 \phi(x)} \sqrt{\nabla^2 \phi(x)}^{\top} = A S_x^{-1} S_x^{-1} A^{\top} = \nabla^2 \phi(x)$ .

For  $x, x' \in \mathcal{X}$ , we have that

$$\sqrt{\nabla^2 \phi(x')} - \sqrt{\nabla^2 \phi(x)} = A(S_{x'} - S_x) 
= -\left(\frac{a_1 a_1^\top (x' - x)}{(a_1^\top x' - b_1)(a_1^\top x - b_1)} \cdots \frac{a_m a_m^\top (x' - x)}{(a_m^\top x' - b_m)(a_m^\top x - b_m)}\right)$$

Therefore,

$$\begin{split} \|\sqrt{\nabla^2\phi(x')} - \sqrt{\nabla^2\phi(x)}\|_{\mathrm{HS}}^2 &= \sum_{i=1}^m \left\| \frac{a_1 a_1^\top (x'-x)}{(a_1^\top x' - b_1)(a_1^\top x - b_1)} \right\|^2 \\ &= \sum_{i=1}^m (x'-x)^\top \frac{a_1 a_1^\top a_1 a_1^\top}{(a_1^\top x' - b_1)^2 (a_1^\top x - b_1)^2} (x'-x) \\ &= (x'-x)^\top \left( \sum_{i=1}^m \frac{a_1 a_1^\top}{(a_1^\top x' - b_1)^2 (a_1^\top x - b_1)^2} \right) (x'-x) \\ &= (x'-x)^\top A S_{x'}^{-2} S_x^{-2} A^\top (x'-x) \\ &= \|v(x,x')\|^2. \end{split}$$

**Modified self-concordance.** The modified self-concordance parameter is

$$\alpha = \sup_{x,x' \in \mathcal{X}} \frac{\|\sqrt{\nabla^2 \phi(x')} - \sqrt{\nabla^2 \phi(x)}\|_{\mathrm{HS}}^2}{\|\nabla \phi(x') - \nabla \phi(x)\|_2^2}$$

$$= \sup_{x,x' \in \mathcal{X}} \frac{\|v(x,x')\|^2}{v(x,x')^\top A^\top A v(x,x')}$$

$$\leq \sup_{v \in \mathbb{R}^d} \frac{\|v\|^2}{v^\top A^\top A v}$$

$$= \max_{i=1,\dots,d} \frac{1}{\sigma_i^2}$$

$$= \frac{1}{\sigma_d^2}.$$

This shows the modified self-concordance parameter can be as large as  $\frac{1}{\sigma_d^2}$ , by choosing appropriate x, x'. For some polyhedra  $\sigma_d \approx 0$ , so  $\alpha \approx 1/\sigma_d^2$  can be arbitrarily large.

**Example in two dimension.** Let d=2, and consider

$$a_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad a_2 = \begin{pmatrix} \sqrt{1 - \epsilon^2} \\ \epsilon \end{pmatrix}$$

for some small  $\epsilon > 0$ , and  $b_1 = b_2 = 0$ . This defines the intersection of two halfspaces:

$$\mathcal{X} = \{x = (x_1, x_2) \colon x_1 \ge 0, \sqrt{1 - \epsilon^2} x_1 + \epsilon x_2 \ge 0\}.$$

The constraint matrix is  $A = \begin{pmatrix} 1 & \sqrt{1-\epsilon^2} \\ 0 & \epsilon \end{pmatrix}$ . We have  $A^{\top}A = \begin{pmatrix} 1 & \sqrt{1-\epsilon^2} \\ \sqrt{1-\epsilon^2} & 1 \end{pmatrix}$  which has eigenvalues  $\sigma_1^2 = 1 + \sqrt{1-\epsilon^2}$  and  $\sigma_2^2 = 1 - \sqrt{1-\epsilon^2}$ . Note that if  $\epsilon$  is small,  $\sigma_1^2 \approx 2$  and  $\sigma_2^2 \approx \epsilon^2/2$ . The corresponding eigenvectors are  $v_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  and  $v_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ .

Let us choose

$$x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad x' = \begin{pmatrix} a \\ b \end{pmatrix}$$

for some constant  $a,b \in \mathbb{R}$ . For simplicity let  $s=\sqrt{1-\epsilon^2}$ . We require  $x' \in \mathcal{X}$ , so  $a \geq 0$  and  $b \geq -\frac{s}{\epsilon}a$ . We have

$$A^{\top}x = \begin{pmatrix} 1 \\ s \end{pmatrix}, \quad A^{\top}x' = \begin{pmatrix} a \\ sa + \epsilon b \end{pmatrix}$$

and

$$A^{\top}(x'-x) = \begin{pmatrix} a-1\\ s(a-1) + \epsilon b \end{pmatrix}.$$

We also have

$$S_x = \begin{pmatrix} 1 & 0 \\ 0 & s \end{pmatrix}, \quad S_x = \begin{pmatrix} a & 0 \\ 0 & sa + \epsilon b \end{pmatrix}.$$

Then

$$\begin{split} v(x,x') &= S_x^{-1} S_{x'}^{-1} A^\top (x'-x) \\ &= \begin{pmatrix} \frac{1}{a} & 0 \\ 0 & \frac{1}{s(sa+\epsilon b)} \end{pmatrix} \begin{pmatrix} a-1 \\ s(a-1)+\epsilon b \end{pmatrix} \\ &= \begin{pmatrix} \frac{a-1}{a} \\ \frac{s(a-1)+\epsilon b}{s(sa+\epsilon b)} \end{pmatrix} \end{split}$$

We want this to be proportional to  $v_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ , so we want

$$\frac{a-1}{a} + \frac{s(a-1) + \epsilon b}{s(sa + \epsilon b)} = 0. \tag{20}$$

We can solve for b in terms of a:

$$b = -\frac{a(a-1)s(s+1)}{\epsilon((a-1)s+a)}.$$

We can see that for all  $a \ge 0$ , this choice of b satisfies the constraint  $b \ge -\frac{s}{\epsilon}a$ , so  $x' \in \mathcal{X}$ . Explicitly, we can choose

$$a = 2$$

$$b = -\frac{2s(s+1)}{\epsilon(s+2)}$$

which satisfies the condition  $b \ge -\frac{2s}{\epsilon}$ . We can verify directly that the condition (20) holds:

$$\frac{a-1}{a} + \frac{s(a-1) + \epsilon b}{s(sa + \epsilon b)} = \frac{1}{2} + \frac{s - \frac{2s(s+1)}{(s+2)}}{s(2s - \frac{2s(s+1)}{(s+2)})} = \frac{1}{2} + \frac{-\frac{s^2}{s+2}}{s(\frac{2s}{(s+2)})} = \frac{1}{2} - \frac{1}{2} = 0.$$

Then with this choice

$$x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad x' = \begin{pmatrix} 2 \\ -\frac{2s(s+1)}{\epsilon(s+2)} \end{pmatrix}$$

we have that  $v(x,x')=\frac{1}{2}v_2$ , i.e. proportional to the eigenvector of  $A^\top A$  with small eigenvalue  $\sigma_2^2$ . Then  $A^\top Av(x,x')=\sigma_2^2v(x,x')$ , and this gives the bound for the modified self-concordance parameter:

$$\alpha \geq \frac{\|v(x,x')\|^2}{v(x,x')^\top A^\top A v(x,x')} = \frac{\|v(x,x')\|^2}{\sigma_2^2 \|v(x,x')\|^2} = \frac{1}{\sigma_2^2} = \frac{1}{1 - \sqrt{1 - \epsilon^2}} \approx \frac{2}{\epsilon^2}.$$

Thus, by setting  $\epsilon \to 0$  we can make  $\alpha$  arbitrarily large. However, note that the case  $\epsilon = 0$  is nice and we have  $\alpha = 1$ , because the domain is a half-space and the problem reduces to one dimension. This example shows the definition of modified self-concordance is not stable.