
Convergence in KL and Rényi Divergence of the Unadjusted Langevin Algorithm Using Estimated Score

Kaylee Yingxi Yang

Department of Statistics and Data Science
Yale University
yingxi.yang@yale.edu

Andre Wibisono

Department of Computer Science
Yale University
andre.wibisono@yale.edu

Abstract

We study Inexact Langevin Algorithm (ILA) for sampling using an estimated score function when the target distribution satisfies log-Sobolev inequality (LSI), motivated by Score-based Generative Modeling (SGM). We prove convergence in Kullback-Leibler (KL) divergence under a sufficient assumption on the error of score estimator called bounded Moment Generating Function (MGF) assumption. Our assumption is weaker than the previous assumption which requires the error has finite L^∞ norm everywhere. Under the L^∞ error assumption, we also prove convergence in Rényi divergence, which is stronger than KL divergence. On the other hand, under L^p error assumption for any $1 \leq p < \infty$ which is weaker than bounded MGF assumption, we show that the stationary distribution of Langevin dynamics with an L^p -accurate score estimator can be far away from the desired distribution. Thus having an L^p -accurate score estimator cannot guarantee convergence. Our results suggest controlling mean squared error which is the form of commonly used loss function when using neural network to estimate score function is not enough to guarantee the upstream algorithm will converge, hence in order to get a theoretical guarantee we need a stronger control over the error in score matching. Despite requiring an exponentially decaying error probability, we give an example to demonstrate the bounded MGF assumption is achievable when using Kernel Density Estimation (KDE)-based score estimator.

1 Introduction

Score-based Generative Modeling (SGM) is a family of sampling methods which have state-of-the-art performance in many applied areas (Song and Ermon, 2019; Ho et al., 2020; Song et al., 2021), but the theoretical understanding of the methods is still lacking. The basic idea of SGM is a two-step procedure. The first step is to estimate the score function from the data (for example via score matching using neural network), the second step is to get new samples from the estimated score via Langevin dynamics. There have been extensive studies on convergence rate of Langevin dynamics and related discrete time algorithms in the setting of knowing exact score function (Vempala and Wibisono, 2019; Chewi et al., 2022). In contrast, theoretical analysis in the setting of using estimated score is limited. De Bortoli et al. (2021) studied the convergence in total variation that requires L^∞ error assumption. Although L^∞ is sufficient to guarantee the convergence, it is too strong to hold in practice since it requires a finite and uniform error at every point x . Lee et al. (2022) and De Bortoli (2022) studied the convergence in total variation and Wasserstein distance of order one respectively under L^2 error assumption. The convergence results are achieved after running the algorithm for a moderate amount of time. The reason why we cannot run the algorithm for infinitely long time is that the stationary distribution of Langevin dynamics with an L^2 -accurate score estimator can be far

away from the desired distribution (Balasubramanian et al., 2022; Lee et al., 2022). This suggests having an L^2 -accurate score estimator is not sufficient to guarantee a distribution which is close to the desired one after running the algorithm for a long time. Therefore there is a need to find a sufficient assumption on the error of score estimator which is achievable in practice and meanwhile can guarantee convergence.

Our contribution We prove convergence in KL divergence under a sufficient assumption on the error of score estimator called bounded Moment Generating Function (MGF) assumption which is weaker than L^∞ but stronger than L^p for $1 \leq p < \infty$. We also prove convergence in Rényi divergence under L^∞ error assumption, which is stronger than KL divergence. Contrary to previous results (Lee et al., 2022; De Bortoli, 2022), our convergence result is stable. Moreover our results hold in KL divergence, which is stronger than TV. We also show that L^p error assumption for any $1 \leq p < \infty$ is not sufficient to guarantee long-term convergence even in TV. Our analysis in this paper suggests controlling mean squared loss in score matching is not enough, we may need a stronger control over the error when constructing the loss function for neural network to train the score function, such as MGF of the error (mean of exponential of squared error). We show the bounded MGF assumption is achievable by giving an example when using Kernel Density Estimator (KDE) to approximate the score function in Gaussian data. Since our original submission to this workshop, we have generalized our analysis to score-based generative models and proved a similar stable convergence result in KL divergence (more details in Wibisono and Yang (2022)).

Notations In this paper, we use $H_\nu(\rho)$ to denote KL divergence of ρ w.r.t. ν , $J_\nu(\rho)$ to denote relative Fisher information of ρ w.r.t. ν , $R_{q,\nu}(\rho)$ to denote Rényi divergence of order q of ρ w.r.t. ν . We provide a review on the definitions in Appendix A.

2 Preliminaries

Given the target distribution $\nu \propto e^{-f}$ in \mathbb{R}^n , the Unadjusted Langevin Algorithm (ULA) with step size $h > 0$ is the following discrete-time algorithm

$$x_{k+1} = x_k + h\nabla \log \nu(x_k) + \sqrt{2h}z_k \quad (1)$$

where $z_k \sim N(0, I)$ are independent standard Gaussians in \mathbb{R}^n . This is a discretization of the Langevin dynamics in continuous time, which is a stochastic process that converges to the target distribution. The convergence properties of the Langevin dynamics have been well studied under various assumptions such as strong log-concavity or weaker isoperimetric inequality such as log-Sobolev inequality (LSI), which allows for some non-log-concavity. There are also convergence guarantees in discrete time with small bias under LSI and smoothness with small step size (Vempala and Wibisono, 2019; Balasubramanian et al., 2022).

In practice, for example in SGM, we may not know true score function $\nabla \log \nu(x)$, so we replace it by an estimated score function $s(x)$ in Eq. (1) and we run Inexact Langevin Algorithm (ILA) with $s(x)$ as follows

$$x_{k+1} = x_k + hs(x_k) + \sqrt{2h}z_k. \quad (2)$$

We study the convergence of above algorithm. One step of the algorithm is the solution $x_{k+1} = X_h$ of the following SDE at time h starting from $X_0 = x_k$

$$dX_t = s(X_t)dt + \sqrt{2}dW_t \quad (3)$$

where W_t is the standard Brownian motion in \mathbb{R}^n . We present the following lemma on Eq. (3), which will be used later.

Lemma 1. *Let ρ_t be the law of Eq. (3), then we have the following bound for time derivative of KL divergence*

$$\frac{d}{dt}H_\nu(\rho_t) \leq -\frac{3}{4}J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(X_0) - \nabla \log \nu(X_t)\|^2 \right].$$

Proof of Lemma 1 is in Appendix B.

3 Main results

3.1 Convergence under bounded MGF assumption

We first introduce the following assumptions.

Assumption 1 (LSI). *The target distribution ν satisfies LSI with constant $\alpha > 0$, which means for any probability distribution ρ*

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

Assumption 2 (L -smoothness). *$f = -\log \nu$ is L -smooth, which means ∇f is L -Lipschitz.*

Assumption 3 (Lipschitz score estimator). *The score estimator $s(x)$ is L_s -Lipschitz.*

Assumption 4 (MGF error assumption). *The error of score estimator $s(x)$ has finite moment generating function of some order r , i.e.*

$$M^2 = \mathbb{E}_\nu \left[\exp \left(r \|\nabla \log \nu(x) - s(x)\|^2 \right) \right] < \infty.$$

Later in this paper we will compare the above assumption with the following.

Assumption 5 (L^∞ error assumption). *The error of score estimator $s(x)$ has finite L^∞ norm at every x , i.e.*

$$M_\infty = \sup_{x \in \mathbb{R}^n} \|\nabla \log \nu(x) - s(x)\| < \infty.$$

Assumption 6 (L^p error assumption). *The error of score estimator $s(x)$ has finite moment of order p for some $1 \leq p < \infty$, i.e.*

$$M_p = \mathbb{E}_\nu \left[\|\nabla \log \nu(x) - s(x)\|^p \right] < \infty.$$

Our main result is the following, which shows biased convergence rate of ILA under the MGF error assumption.

Theorem 1 (Convergence of KL divergence under MGF error assumption). *Suppose Assumption 1-3 and 4 with $r = \frac{9}{\alpha}$ hold. If $0 < h < \min(\frac{9}{12L_sL}, \frac{1}{2\alpha})$, then after k iterations of ILA (2)*

$$H_\nu(\rho_k) \leq e^{-\frac{1}{4}\alpha h k} H_\nu(\rho_0) + C_1 h + C_2 \log M$$

where $C_1 = O(\frac{nL_s^2L}{\alpha})$ and $C_2 = \frac{16}{3}$.

Proof of Theorem 1 is in Appendix C. The convergence rate is similar to Theorem 1 in Vempala and Wibisono (2019) in the setting of knowing exact score function but has an extra non-vanishing term induced by the error of score estimator. So in order to have a small asymptotic error, we need an accurate score estimator. Theorem 1 directly implies the following corollary.

Corollary 1. *Suppose Assumption 1-3 hold. For any $\varepsilon > 0$, if the score estimator $s(x)$ has a small MGF error*

$$\log \mathbb{E}_\nu \left[\exp \left(\frac{9}{\alpha} \|\nabla \log \nu(x) - s(x)\|^2 \right) \right] \leq \frac{3}{16} \varepsilon,$$

then running ILA (2) with step size $h \leq \frac{\varepsilon}{4C_1} = O\left(\frac{\varepsilon\alpha}{nL_s^2L}\right)$ for at least $k = \frac{4}{\alpha h} \log \frac{4H_\nu(\rho_0)}{\varepsilon}$ iterations reaches $H_\nu(\rho_k) \leq \varepsilon$.

3.2 Comparing different error assumptions

Because L^∞ implies MGF error assumption, the convergence result in Theorem 1 also holds under Assumption 5 in place of Assumption 4. In this case, we provide a simpler proof in Appendix D. In addition, under L^∞ error assumption, we prove convergence in Rényi divergence of order $q \geq 1$, which is stronger than KL divergence. We conjecture the convergence holds under the MGF error assumption and leave it for future work.

Theorem 2 (Convergence of Rényi divergence under L^∞ error assumption). *Suppose Assumption 1-3 and 5 hold. Let $q \geq 1$. If $0 < h < \min(\frac{\alpha}{12L_sLq}, \frac{q}{4\alpha})$, then after k iterations of ILA (2)*

$$R_{q,\nu}(\rho_k) \leq e^{-\frac{q}{2}h k} R_{q,\nu}(\rho_0) + C_1 h + C_2 M_\infty^2,$$

where $C_1 = O(\frac{hL_s^2Lq^2}{\alpha})$ and $C_2 = O(\frac{q^2}{\alpha})$.

Proof of Theorem 2 is in Appendix E.

L^∞ -accurate score estimator requires a finite and uniform error at every point x , which is too strong and may not hold in practice. By the construction of loss function in score estimation, the estimated score is only L^2 -accurate with high probability (Proposition 9 in Block et al. (2020)). However, the stationary distribution of Langevin dynamics with L^2 -accurate score can be far away from the true distribution in TV (Lee et al., 2022; Balasubramanian et al., 2022). So having an L^2 -accurate score estimator cannot guarantee the sampling algorithm will converge to the desired distribution. We further claim that any finite higher moment (Assumption 6) cannot guarantee the convergence either. This is to say, the stationary distribution of Langevin dynamics with an L^p -accurate score estimator where $1 \leq p < \infty$ can also be far away from the true distribution. We use the example in Balasubramanian et al. (2022) to illustrate this.

Example 1 (L^p bound is not sufficient). Let $\nu = \frac{3}{4}\mathcal{N}(-m, 1) + \frac{1}{4}\mathcal{N}(m, 1)$ and $\mu = \frac{1}{2}\mathcal{N}(-m, 1) + \frac{1}{2}\mathcal{N}(m, 1)$. For all $m \geq \frac{1}{80}$ and $p \geq 1$, the following holds

$$\mathbb{E}_\nu \left[\|\nabla \log \nu - \nabla \log \mu\|^p \right] \leq 4^{p-1} m^p \exp\left(-\frac{m^2}{2}\right) \rightarrow 0 \text{ as } m \rightarrow \infty.$$

However, the total variation between μ and ν is large, $TV(\mu, \nu) \geq \frac{1}{800}$.

The example above shows if the target distribution is ν , but we run Langevin dynamics with an L^p -accurate score estimator $s = \nabla \log \mu$, then the limiting distribution will be μ , hence $TV(\mu, \nu)$ is the asymptotic bias. We provide more details of Example 1 in Appendix F. Therefore, in order to obtain a convergence guarantee with an achievable sufficient assumption, we need the error of score estimator in between L^p and L^∞ .

The MGF error assumption is weaker than L^∞ but stronger than L^p error bound for any $1 \leq p < \infty$. Indeed, the MGF error assumption implies $\mathbb{P}_\nu(\|\nabla \log \nu(x) - s(x)\|^2 \geq c_n) \rightarrow 0$ exponentially as $c_n \rightarrow \infty$, which is strong. It is not clear whether the assumption holds via score matching using neural network. But here we show an example where the assumption is achieved when using a simple KDE-based score estimator in Gaussian data.

Example 2. Let $\nu = N(\mu, \sigma^2)$, then $f = -\log \nu$ is $1/\sigma^2$ -strongly convex, thus ν satisfies LSI with constant $1/\sigma^2$. The score function is $\nabla \log \nu(x) = -(x - \mu)/\sigma^2$. We estimate the score function via Gaussian KDE with bandwidth $h > 0$, so the score estimator is

$$s(x) = \nabla \log(\nu * N(0, h^2)) = -\frac{x - \mu}{\sigma^2 + h^2}$$

and $s(x)$ is $\frac{1}{\sigma^2 + h^2}$ -Lipschitz. Let $r = 9\sigma^2$. If $h < \frac{\sigma}{\sqrt{3\sqrt{2}-1}}$, then the MGF is finite,

$$\mathbb{E}_\nu \left[\exp \left(r \|\nabla \log \nu(x) - s(x)\|^2 \right) \right] = \mathbb{E}_{N(0,1)} \left[\exp \left(\frac{9h^4}{(h^2 + \sigma^2)^2} Z^2 \right) \right] = \frac{h^2 + \sigma^2}{\sqrt{\sigma^4 + 2h^2\sigma^2 - 17h^4}}.$$

4 Discussion

In this paper we study the convergence of Inexact Langevin Algorithm (ILA) using estimated score function, motivated by SGM. We prove convergence result in KL divergence under MGF error assumption, which is weaker than L^∞ error assumption. In addition, under L^∞ assumption, we provide convergence in Rényi divergence, which is stronger than KL divergence. On the other hand, we show L^p error assumption for any $1 \leq p < \infty$ is not sufficient to guarantee convergence. Our results suggest controlling L^2 loss (mean squared error) in score matching is not enough to guarantee the upstream algorithm will converge, therefore in order to get a theoretical guarantee one may need a stronger control over the error when constructing the loss function for neural network to estimate the score function, such as MGF (mean of exponential of squared error). Using such loss functions will give us a more accurate score estimator especially in low density regions, since they punish more on those large errors which are usually from low density areas (Song and Ermon, 2019). Although the MGF error assumption requires an exponentially decaying tail, we show it is achievable when using a simple KDE score estimator in Gaussian data. It is intriguing to explore if the KDE score estimator satisfies the assumption for a larger class of distributions, such as strongly log-concave distributions. It is also interesting to investigate if the score matching estimator by neural network satisfies this assumption.

References

- Krishna Balasubramanian, Sinho Chewi, Murat A Erdogdu, Adil Salim, and Shunshi Zhang. Towards a theory of non-log-concave sampling: First-order stationarity guarantees for Langevin Monte Carlo. In *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 2896–2923, 02–05 Jul 2022.
- Adam Block, Youssef Mroueh, and Alexander Rakhlin. Generative modeling with denoising auto-encoders and Langevin sampling. *arXiv preprint arXiv:2002.00107*, 2020.
- Sinho Chewi, Murat A. Erdogdu, Mufan Li, Ruoqi Shen, and Shunshi Zhang. Analysis of Langevin Monte Carlo from Poincare to log-Sobolev. In *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 1–2, 02–05 Jul 2022.
- Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. *arXiv preprint arXiv:2208.05314*, 2022.
- Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion Schrödinger Bridge with applications to score-based generative modeling. In *Advances in Neural Information Processing Systems*, 2021.
- Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The Total Variation distance between high-dimensional Gaussians. *arXiv preprint arXiv:1810.08693*, 2018.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851, 2020.
- Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with polynomial complexity. *arXiv preprint arXiv:2206.06227*, 2022.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- Santosh Vempala and Andre Wibisono. Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry suffices. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Andre Wibisono and Kaylee Yingxi Yang. Convergence in KL divergence of the inexact Langevin algorithm with application to score-based generative models. *arXiv preprint arXiv:2211.01512*, 2022.

A A review on notations and definitions

In this section, we review notations and definitions of Kullback-Leibler (KL) divergence, relative Fisher information, Rényi divergence and Rényi information. Let ρ, ν be two probability distributions in \mathbb{R}^n denoted by their probability density functions w.r.t. Lebesgue measure on \mathbb{R}^n . Assume ρ and ν have full support on \mathbb{R}^n .

Definition 1 (KL divergence). *The Kullback-Leibler (KL) divergence of ρ w.r.t. ν is*

$$H_\nu(\rho) = \int_{\mathbb{R}^n} \rho \log \frac{\rho}{\nu} dx.$$

Definition 2 (Relative Fisher information). *The relative Fisher information of ρ w.r.t. ν is*

$$J_\nu(\rho) = \int_{\mathbb{R}^n} \rho \|\nabla \log \frac{\rho}{\nu}\|^2 dx.$$

Definition 3 (Rényi divergence). *For $q \geq 0, q \neq 1$, the Rényi divergence of order q of ρ w.r.t. ν is*

$$R_{q,\nu}(\rho) = \frac{1}{q-1} \log F_{q,\nu}(\rho)$$

where

$$F_{q,\nu}(\rho) = \mathbb{E}_\nu \left[\left(\frac{\rho}{\nu} \right)^q \right].$$

Definition 4 (Rényi information). *For $q \geq 0$, the Rényi information of order q of ρ w.r.t. ν is*

$$G_{q,\nu}(\rho) = \mathbb{E}_\nu \left[\left(\frac{\rho}{\nu} \right)^q \|\nabla \log \frac{\rho}{\nu}\|^2 \right] = \frac{4}{q^2} \mathbb{E}_\nu \left[\|\nabla \left(\frac{\rho}{\nu} \right)^{\frac{q}{2}}\|^2 \right].$$

B Proof of Lemma 1

Proof. The continuity equation corresponding to Eq. (3) is

$$\frac{\partial \rho_t(x)}{\partial t} = -\nabla \cdot \left(\rho_t(x) \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] \right) + \Delta \rho_t(x).$$

Therefore

$$\begin{aligned} \frac{d}{dt} H_\nu(\rho_t) &= \int_{\mathbb{R}^n} \frac{\partial \rho_t}{\partial t} \log \frac{\rho_t}{\nu} dx \\ &= \int_{\mathbb{R}^n} \left(-\nabla \cdot \left(\rho_t \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] \right) + \Delta \rho_t \right) \log \frac{\rho_t}{\nu} dx \\ &= \int_{\mathbb{R}^n} \left(-\nabla \cdot \left(\rho_t \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] \right) + \nabla \cdot \left(\rho_t \nabla \log \frac{\rho_t}{\nu} \right) + \nabla \cdot \left(\rho_t \nabla \log \nu \right) \right) \log \frac{\rho_t}{\nu} dx \\ &= \int_{\mathbb{R}^n} \left(\nabla \cdot \left(\rho_t \left(\nabla \log \frac{\rho_t}{\nu} - \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] + \nabla \log \nu \right) \right) \right) \log \frac{\rho_t}{\nu} dx \\ &= - \int_{\mathbb{R}^n} \rho_t \langle \nabla \log \frac{\rho_t}{\nu} - \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] + \nabla \log \nu, \nabla \log \frac{\rho_t}{\nu} \rangle dx \quad \text{integrating by parts} \\ &= - \int_{\mathbb{R}^n} \rho_t \|\nabla \log \frac{\rho_t}{\nu}\|^2 dx + \int_{\mathbb{R}^n} \rho_t \langle \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] - \nabla \log \nu, \nabla \log \frac{\rho_t}{\nu} \rangle dx \\ &= -J_\nu(\rho_t) + \int_{\mathbb{R}^n} \rho_t \langle \mathbb{E}_{\rho_{0|t}}[s(x_0)|x_t = x] - \nabla \log \nu, \nabla \log \frac{\rho_t}{\nu} \rangle dx \\ &= -J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\langle s(x_0) - \nabla \log \nu(x_t), \nabla \log \frac{\rho_t(x_t)}{\nu(x_t)} \rangle \right] \quad \text{by renaming } x \text{ as } x_t \\ &\leq -J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - \nabla \log \nu(x_t)\|^2 \right] + \frac{1}{4} \mathbb{E}_{\rho_{0t}} \left[\|\nabla \log \frac{\rho_t(x_t)}{\nu(x_t)}\|^2 \right] \quad \text{by } \langle a, b \rangle \leq \|a\|^2 + \frac{1}{4} \|b\|^2 \\ &= -J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - \nabla \log \nu(x_t)\|^2 \right] + \frac{1}{4} J_\nu(\rho_t) \\ &= -\frac{3}{4} J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - \nabla \log \nu(x_t)\|^2 \right]. \end{aligned}$$

□

C Proof of Theorem 1

We restate the full theorem here.

Theorem 1. *Suppose the target distribution $\nu \propto e^{-f}$ satisfies LSI with constant $\alpha > 0$ and f is L -smooth. Let $s(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be an L_s -Lipschitz score estimator and the error satisfies Assumption 4. Let $0 < h < \min(\frac{\alpha}{12L_sL}, \frac{1}{2\alpha})$, after k iterations of ILA (2)*

$$H_\nu(\rho_k) \leq e^{-\frac{1}{4}\alpha h k} H_\nu(\rho_0) + \frac{768nL_s^2L}{\alpha} h^2 + \frac{128nL_s^2}{\alpha} h + \frac{16}{3} \log M.$$

To prove Theorem 1, we will use the following auxiliary result.

Lemma 2. *If the score estimator $s(x)$ is L_s -Lipschitz and $t \leq \frac{1}{3L_s}$, then*

$$\|s(x_t) - s(x_0)\|^2 \leq 18L_s^2t^2\|s(x_t) - \nabla \log \nu(x_t)\|^2 + 18L_s^2t^2\|\nabla \log \nu(x_t)\|^2 + 6L_s^2t\|z_0\|^2.$$

Lemma 3. *Suppose the assumptions in Theorem 1 hold, then along each step of ILA (2),*

$$H_\nu(\rho_{k+1}) \leq e^{-\frac{1}{4}\alpha h} H_\nu(\rho_k) + 144nL_s^2Lh^3 + 24nL_s^2h^2 + \alpha h \log M.$$

Proof of Lemma 2. By L_s -Lipschitzness

$$\|s(x_t) - s(x_0)\|^2 \leq L_s^2\|x_t - x_0\|^2 = L_s^2\|ts(x_0) + \sqrt{2t}z_0\|^2 \leq 2L_s^2t^2\|s(x_0)\|^2 + 4L_s^2t\|z_0\|^2.$$

It is more convenient for our subsequent analysis to have a bound in terms of $s(x_t)$ rather than $s(x_0)$, so we use

$$L_s\|x_t - x_0\| \geq \|s(x_t) - s(x_0)\| \geq \|s(x_0)\| - \|s(x_t)\|.$$

Rearranging it gives

$$\begin{aligned} \|s(x_0)\| &\leq L_s\|x_t - x_0\| + \|s(x_t)\| \\ &= L_s\|ts(x_0) + \sqrt{2t}z_0\| + \|s(x_t)\| \quad \text{since } x_t = x_0 + ts(x_0) + \sqrt{2t}z_0 \\ &= L_s t \|s(x_0)\| + L_s \sqrt{2t} \|z_0\| + \|s(x_t)\| \quad \text{by triangle inequality} \\ &\leq \frac{1}{3} \|s(x_0)\| + L_s \sqrt{2t} \|z_0\| + \|s(x_t)\| \quad \text{since } t \leq \frac{1}{3L_s} \iff t \leq h \leq \frac{\alpha}{12L_sL} \text{ and } \alpha \leq L \end{aligned}$$

Therefore

$$\begin{aligned} \|s(x_0)\| &\leq \frac{3}{2} \|s(x_t)\| + \frac{3}{\sqrt{2}} L_s \sqrt{t} \|z_0\| \\ \implies \|s(x_0)\|^2 &\leq \frac{9}{2} \|s(x_t)\|^2 + 9L_s^2t\|z_0\|^2 \quad \text{by } \|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2 \end{aligned} \tag{4}$$

So we can bound $\|s(x_t) - s(x_0)\|^2$ as follows

$$\begin{aligned} \|s(x_t) - s(x_0)\|^2 &\leq 2L_s^2t^2\|s(x_0)\|^2 + 4L_s^2t\|z_0\|^2 \\ &\leq 2L_s^2t^2 \left(\frac{9}{2} \|s(x_t)\|^2 + 9L_s^2t\|z_0\|^2 \right) + 4L_s^2t\|z_0\|^2 \quad \text{by plugging in Eq. (4)} \\ &= 9L_s^2t^2\|s(x_t)\|^2 + (18L_s^4t^3 + 4L_s^2t)\|z_0\|^2 \\ &\leq 9L_s^2t^2\|s(x_t)\|^2 + 6L_s^2t\|z_0\|^2 \quad \text{since } t \leq \frac{1}{3L_s} \\ &= 9L_s^2t^2\|s(x_t) - \nabla \log \nu(x_t) + \nabla \log \nu(x_t)\|^2 + 6L_s^2t\|z_0\|^2 \\ &\leq 18L_s^2t^2\|s(x_t) - \nabla \log \nu(x_t)\|^2 + 18L_s^2t^2\|\nabla \log \nu(x_t)\|^2 + 6L_s^2t\|z_0\|^2. \end{aligned}$$

□

Proof of Lemma 3. For simplicity suppose $k = 0$, then one step of ULA using estimated score with step size $h > 0$ is

$$x_1 \stackrel{d}{=} x_0 + hs(x_0) + \sqrt{2h}z_0.$$

This is the solution to the following SDE at time $t = h$

$$dX_t = s(X_t)dt + \sqrt{2}dW_t$$

where W_t is the standard Brownian motion in \mathbb{R}^n .

Let $M(x) = \|\nabla \log \nu(x) - s(x)\|^2$. By Lemma 1,

$$\begin{aligned} \frac{\partial}{\partial t} H_\nu(\rho_t) &\leq -\frac{3}{4}J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - \nabla \log \nu(x_t)\|^2 \right] \\ &\leq -\frac{3}{4}J_\nu(\rho_t) + 2\mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - s(x_t)\|^2 \right] + 2\mathbb{E}_{\rho_t} \left[\|s(x_t) - \nabla \log \nu(x_t)\|^2 \right] \\ &\leq -\frac{3}{4}J_\nu(\rho_t) + 2\mathbb{E}_{\rho_{0t}} \left[18L_s^2 t^2 \|s(x_t) - \nabla \log \nu(x_t)\|^2 + 18L_s^2 t^2 \|\nabla \log \nu(x_t)\|^2 + 6L_s^2 t \|z_0\|^2 \right] \\ &\quad + 2\mathbb{E}_{\rho_t} \left[M(x) \right] \quad \text{by Lemma 2} \\ &= -\frac{3}{4}J_\nu(\rho_t) + (36L_s^2 t^2 + 2)\mathbb{E}_{\rho_t} \left[M(x) \right] + 36L_s^2 t^2 \mathbb{E}_{\rho_t} \left[\|\nabla \log \nu(x_t)\|^2 \right] + 12nL_s^2 t \\ &\stackrel{(i)}{\leq} -\frac{3}{4}J_\nu(\rho_t) + \frac{9}{4}\mathbb{E}_{\rho_t} \left[M(x) \right] + 36L_s^2 t^2 \mathbb{E}_{\rho_t} \left[\|\nabla \log \nu(x_t)\|^2 \right] + 12nL_s^2 t \\ &\stackrel{(ii)}{\leq} -\frac{3}{4}J_\nu(\rho_t) + \frac{9}{4}\mathbb{E}_{\rho_t} \left[M(x) \right] + 36L_s^2 t^2 \left(\frac{4L^2}{\alpha} H_\nu(\rho_t) + 2nL \right) + 12nL_s^2 t \\ &= -\frac{3}{4}J_\nu(\rho_t) + \frac{9}{4}\mathbb{E}_{\rho_t} \left[M(x) \right] + \frac{144L_s^2 t^2 L^2}{\alpha} H_\nu(\rho_t) + 72nL_s^2 t^2 L + 12nL_s^2 t \\ &\leq -\frac{3}{4}J_\nu(\rho_t) + \frac{9}{4}\mathbb{E}_{\rho_t} \left[M(x) \right] + \alpha H_\nu(\rho_t) + 72nL_s^2 t^2 L + 12nL_s^2 t \quad \text{since } t^2 \leq h^2 \leq \frac{\alpha^2}{144L_s^2 L^2} \\ &\leq -\frac{3}{2}\alpha H_\nu(\rho_t) + \frac{9}{4}\mathbb{E}_{\rho_t} \left[M(x) \right] + \alpha H_\nu(\rho_t) + 72nL_s^2 t^2 L + 12nL_s^2 t \quad \text{by LSI} \\ &= -\frac{1}{2}\alpha H_\nu(\rho_t) + \frac{9}{4}\mathbb{E}_{\rho_t} \left[M(x) \right] + 72nL_s^2 t^2 L + 12nL_s^2 t. \end{aligned}$$

where (i) is because $t^2 \leq h^2 \leq \frac{\alpha^2}{144L_s^2 L^2} \leq \frac{1}{144L_s^2}$ and (ii) is by Lemma 12 in Vempala and Wibisono (2019). We then bound the second term by applying Donsker–Varadhan variational characterizations of KL divergence $\mathbb{E}_P[f(x)] \leq \log \mathbb{E}_Q e^{f(x)} + H_Q(P)$ for $f(x) = \frac{9}{\alpha}M(x)$, $P = \rho_t$ and $Q = \nu$

$$\begin{aligned} \mathbb{E}_{\rho_t} \left[\frac{9}{\alpha}M(x) \right] &\leq \log \mathbb{E}_\nu \left[e^{\frac{9}{\alpha}M(x)} \right] + H_\nu(\rho_t) \\ \iff \mathbb{E}_{\rho_t} \left[M(x) \right] &\leq \frac{\alpha}{9} \log \mathbb{E}_\nu \left[e^{\frac{9}{\alpha}M(x)} \right] + \frac{\alpha}{9} H_\nu(\rho_t). \end{aligned}$$

Therefore

$$\begin{aligned} \frac{\partial}{\partial t} H_\nu(\rho_t) &\leq -\frac{1}{2}\alpha H_\nu(\rho_t) + \frac{\alpha}{4} \log \mathbb{E}_\nu \left[e^{\frac{9}{\alpha}M(x)} \right] + \frac{1}{4}\alpha H_\nu(\rho_t) + 72nL_s^2 t^2 L + 12nL_s^2 t \\ &= -\frac{1}{4}\alpha H_\nu(\rho_t) + \frac{\alpha}{4} \log \mathbb{E}_\nu \left[e^{\frac{9}{\alpha}M(x)} \right] + 72nL_s^2 t^2 L + 12nL_s^2 t \\ &\leq -\frac{1}{4}\alpha H_\nu(\rho_t) + \frac{1}{2}\alpha \log M + 72nL_s^2 t^2 L + 12nL_s^2 t \quad \text{by MGF error assumption} \\ &\leq -\frac{1}{4}\alpha H_\nu(\rho_t) + \frac{1}{2}\alpha \log M + 72nL_s^2 h^2 L + 12nL_s^2 h \quad \text{since } t \in (0, h) \end{aligned}$$

This is equivalent to

$$\begin{aligned} \frac{\partial}{\partial t} e^{\frac{1}{4}\alpha t} H_\nu(\rho_t) &\leq e^{\frac{1}{4}\alpha t} \left(\frac{1}{2}\alpha \log M + 72nL_s^2 h^2 L + 12nL_s^2 h \right) \\ \implies e^{\frac{1}{4}\alpha h} H_\nu(\rho_h) &\leq H_\nu(\rho_0) + \frac{4(e^{\frac{1}{4}\alpha h} - 1)}{\alpha} \left(\frac{1}{2}\alpha \log M + 72nL_s^2 h^2 L + 12nL_s^2 h \right) \\ \implies H_\nu(\rho_h) &\leq e^{-\frac{1}{4}\alpha h} H_\nu(\rho_0) + 2h \left(\frac{1}{2}\alpha \log M + 72nL_s^2 h^2 L + 12nL_s^2 h \right) \end{aligned}$$

since $e^{-\frac{1}{4}\alpha h} \leq 1$ and $e^c - 1 \leq 2c$ for $c = \frac{\alpha}{4}h \in (0, 1)$. Renaming ρ_0 as ρ_k and ρ_h as ρ_{k+1} ,

$$H_\nu(\rho_{k+1}) \leq e^{-\frac{1}{4}\alpha h} H_\nu(\rho_k) + 144nL_s^2Lh^3 + 24nL_s^2h^2 + \alpha h \log M.$$

□

The proof of Theorem 1 then follows Lemma 3.

Proof of Theorem 1. Applying the recursion contraction k times

$$\begin{aligned} H_\nu(\rho_k) &\leq e^{-\frac{1}{4}\alpha hk} H_\nu(\rho_0) + \sum_{i=0}^{k-1} e^{-\frac{1}{4}\alpha hi} \left(\alpha h \log M + 144nL_s^2h^3L + 24nL_s^2h^2 \right) \\ &\leq e^{-\frac{1}{4}\alpha hk} H_\nu(\rho_0) + \frac{1}{1 - e^{-\frac{1}{4}\alpha h}} \left(\alpha h \log M + 144nL_s^2h^3L + 24nL_s^2h^2 \right) \\ &\leq e^{-\frac{1}{4}\alpha hk} H_\nu(\rho_0) + \frac{16}{3\alpha h} \left(\alpha h \log M + 144nL_s^2h^3L + 24nL_s^2h^2 \right) \\ &= e^{-\frac{1}{4}\alpha hk} H_\nu(\rho_0) + \frac{768nL_s^2L}{\alpha} h^2 + \frac{128nL_s^2}{\alpha} h + \frac{16}{3} \log M. \end{aligned}$$

□

D Proof of convergence of KL divergence under L^∞ error assumption

We state the theorem here.

Theorem 3. Suppose the target distribution $\nu \propto e^{-f}$ satisfies LSI with constant $\alpha > 0$ and f is L -smooth. Let $s(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be an L_s -Lipschitz score estimator with error bound in L^∞ , i.e.

$$M_\infty = \max_{x \in \mathbb{R}^n} \|\nabla \log \nu(x) - s(x)\| < \infty.$$

Let $0 < h < \min(\frac{\alpha}{12L_sL}, \frac{1}{2\alpha})$, then after k iterations of ILA (2)

$$H_\nu(\rho_k) \leq e^{-\frac{\alpha}{2}hk} H_\nu(\rho_0) + C_1h + C_2M_\infty^2,$$

for some constants C_1 and C_2 .

We will use the following auxiliary result.

Lemma 4. Suppose the assumptions in Theorem 3 hold, then along each step of ILA (2)

$$H_\nu(\rho_{k+1}) \leq e^{-\frac{\alpha h}{2}} H_\nu(\rho_k) + 144nLL_s^2h^3 + 24nL_s^2h^2 + (72L_s^2h^3 + 4h)M_\infty^2.$$

Proof of Lemma 4.

$$\begin{aligned} \frac{d}{dt} H_\nu(\rho_t) &\leq -\frac{3}{4} J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - \nabla \log \nu(x_t)\|^2 \right] \quad \text{by Lemma 1} \\ &= -\frac{3}{4} J_\nu(\rho_t) + \mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - s(x_t) + s(x_t) - \nabla \log \nu(x_t)\|^2 \right] \\ &\leq -\frac{3}{4} J_\nu(\rho_t) + 2\mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - s(x_t)\|^2 \right] + 2\mathbb{E}_{\rho_t} \left[\|s(x_t) - \nabla \log \nu(x_t)\|^2 \right] \\ &\leq -\frac{3}{4} J_\nu(\rho_t) + 2\mathbb{E}_{\rho_{0t}} \left[\|s(x_0) - s(x_t)\|^2 \right] + 2M_\infty^2 \quad \text{by } L^\infty \text{ error assumption.} \end{aligned}$$

By Lemma 2,

$$\begin{aligned} \|s(x_t) - s(x_0)\|^2 &\leq 18L_s^2t^2\|s(x_t) - \nabla \log \nu(x_t)\|^2 + 18L_s^2t^2\|\nabla \log \nu(x_t)\|^2 + 6L_s^2t\|z_0\|^2 \\ &\leq 18L_s^2t^2M_\infty^2 + 18L_s^2t^2\|\nabla \log \nu(x_t)\|^2 + 6L_s^2t\|z_0\|^2 \quad \text{by } L^\infty \text{ error assumption.} \end{aligned}$$

The condition $t \leq \frac{1}{3L_s}$ in Lemma 2 is satisfied since $t \leq h \leq \frac{\alpha}{12L_sL}$ and $\alpha \leq L$. Therefore,

$$\begin{aligned} 2\mathbb{E}_{\rho_{0t}} \left[\|s(x_t) - s(x_0)\|^2 \right] &\leq 36L_s^2t^2M_\infty^2 + 36L_s^2t^2\mathbb{E}_{\rho_t} \left[\|\nabla \log \nu(x_t)\|^2 \right] + 12nL_s^2t \\ &\leq 36L_s^2t^2M_\infty^2 + 36L_s^2t^2 \left(\frac{4L^2}{\alpha} H_\nu(\rho_t) + 2nL \right) + 12nL_s^2t \\ &= \frac{144L_s^2t^2L^2}{\alpha} H_\nu(\rho_t) + 72nLL_s^2t^2 + 12nL_s^2t + 36L_s^2t^2M_\infty^2. \end{aligned}$$

The second inequality is by Lemma 12 in Vempala and Wibisono (2019). Therefore, the time derivative of KL divergence is bounded by

$$\begin{aligned} \frac{d}{dt} H_\nu(\rho_t) &\leq -\frac{3}{4} J_\nu(\rho_t) + \frac{144L_s^2t^2L^2}{\alpha} H_\nu(\rho_t) + 72nLL_s^2t^2 + 12nL_s^2t + 36L_s^2t^2M_\infty^2 + 2M_\infty^2 \\ &\leq -\frac{3\alpha}{2} H_\nu(\rho_t) + \frac{144L_s^2t^2L^2}{\alpha} H_\nu(\rho_t) + 72nLL_s^2t^2 + 12nL_s^2t + 36L_s^2t^2M_\infty^2 + 2M_\infty^2 && \text{by LSI} \\ &\leq -\frac{3\alpha}{2} H_\nu(\rho_t) + \frac{144L_s^2h^2L^2}{\alpha} H_\nu(\rho_t) + 72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 && \text{since } t \in (0, h) \\ &\leq -\frac{3\alpha}{2} H_\nu(\rho_t) + \alpha H_\nu(\rho_t) + 72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 && \text{since } h^2 \leq \frac{\alpha^2}{144L_s^2L^2} \\ &= -\frac{\alpha}{2} H_\nu(\rho_t) + 72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 \end{aligned}$$

This is equivalent to

$$\frac{d}{dt} e^{\frac{\alpha}{2}t} H_\nu(\rho_t) \leq e^{\frac{\alpha}{2}t} \left(72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 \right).$$

Integrating from 0 to h

$$\begin{aligned} e^{\frac{\alpha}{2}h} H_\nu(\rho_h) - H_\nu(\rho_0) &\leq \frac{2(e^{\frac{\alpha}{2}h} - 1)}{\alpha} \left(72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 \right) \\ &\leq 2h \left(72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 \right). \end{aligned}$$

where the last inequality is from $e^c - 1 \leq 2c$ for $c = \frac{\alpha}{2}h \in (0, 1)$. Rearranging gives

$$\begin{aligned} H_\nu(\rho_h) &\leq e^{-\frac{\alpha}{2}h} H_\nu(\rho_0) + e^{-\frac{\alpha}{2}h} 2h \left(72nLL_s^2h^2 + 12nL_s^2h + 36L_s^2h^2M_\infty^2 + 2M_\infty^2 \right) \\ &\leq e^{-\frac{\alpha}{2}h} H_\nu(\rho_0) + 144nLL_s^2h^3 + 24nL_s^2h^2 + (72L_s^2h^3 + 4h)M_\infty^2 && \text{since } e^{-\frac{\alpha}{2}h} \leq 1. \end{aligned}$$

Rename $\rho_0 \equiv \rho_k, \rho_h \equiv \rho_{k+1}$,

$$H_\nu(\rho_{k+1}) \leq e^{-\frac{\alpha}{2}h} H_\nu(\rho_k) + 144nLL_s^2h^3 + 24nL_s^2h^2 + (72L_s^2h^3 + 4h)M_\infty^2. \quad (5)$$

□

Proof of Theorem 3. Applying the recursion (5) k times, we obtain

$$\begin{aligned} H_\nu(\rho_k) &\leq e^{-\frac{\alpha}{2}hk} H_\nu(\rho_0) + \sum_{i=0}^{k-1} e^{-\frac{\alpha}{2}hi} \left(144nLL_s^2h^3 + 24nL_s^2h^2 + 72L_s^2h^3M_\infty^2 + 4M_\infty^2h \right) \\ &\leq e^{-\frac{\alpha}{2}hk} H_\nu(\rho_0) + \frac{1}{1 - e^{-\frac{\alpha}{2}h}} \left(144nLL_s^2h^3 + 24nL_s^2h^2 + 72L_s^2h^3M_\infty^2 + 4M_\infty^2h \right) \\ &\stackrel{(i)}{\leq} e^{-\frac{\alpha}{2}hk} H_\nu(\rho_0) + \frac{8}{3\alpha h} \left(144nLL_s^2h^3 + 24nL_s^2h^2 + 72L_s^2h^3M_\infty^2 + 4M_\infty^2h \right) \\ &\leq e^{-\frac{\alpha}{2}hk} H_\nu(\rho_0) + \frac{384nLL_s^2}{\alpha} h^2 + \frac{64nL_s^2}{\alpha} h + \left(\frac{192L_s^2h^2}{\alpha} + \frac{32}{3\alpha} \right) M_\infty^2. \end{aligned}$$

where (i) is because $1 - e^{-c} \geq \frac{3}{4}c$ for $0 < c = \frac{\alpha}{2}h < \frac{1}{4}$. □

E Proof of Theorem 2

We restate the full theorem here.

Theorem 2. *Suppose the assumptions in Theorem 3 hold and let $0 < h < \min(\frac{\alpha}{12LL_sq}, \frac{q}{4\alpha})$, then after k iterations of ILA (2)*

$$\begin{aligned} R_{q,\nu}(\rho_k) &\leq e^{-\frac{\alpha}{2}hk} R_{q,\nu}(\rho_0) + \frac{192nLL_s^2q^2}{\alpha}h^2 + \frac{32nL_s^2q^2}{\alpha}h + \left(\frac{96L_s^2h^2q^2}{\alpha} + \frac{16q^2}{3\alpha}\right)M_\infty^2 \\ &\leq e^{-\frac{\alpha}{2}hk} R_{q,\nu}(\rho_0) + C_1h + C_2M_\infty^2, \end{aligned}$$

for some constants C_1 and C_2 .

We will use the following auxiliary result.

Lemma 5. *Let $\varphi_t = \frac{\rho_t}{\nu}$ and $\psi_t = \frac{\varphi_t^{q-1}}{\mathbb{E}_\nu[\varphi_t^q]} = \frac{\varphi_t^{q-1}}{F_{q,\nu}(\rho_t)}$, then*

$$\frac{\partial}{\partial t} R_{q,\nu}(\rho_t) \leq -\frac{3}{4}q \frac{G_{q,\nu}(\rho_t)}{F_{q,\nu}(\rho_t)} + q\mathbb{E}_{\rho_{0t}}[\psi_t(x_t)\|s(x_0) - \nabla \log \nu(x_t)\|^2].$$

This is a generalized version of Proposition 15 in Chewi et al. (2022) to the setting of estimated score.

Lemma 6. *Suppose the assumptions in Theorem 2 hold, then along each step of ILA (2)*

$$R_{q,\nu}(\rho_{k+1}) \leq e^{-\frac{\alpha}{q}h} R_{q,\nu}(\rho_k) + 144nLL_s^2qh^3 + 24nL_s^2qh^2 + (72L_s^2qh^3 + 4qh)M_\infty^2.$$

Proof of Lemma 5.

$$\begin{aligned} \frac{\partial}{\partial t} R_{q,\nu}(\rho_t) &= \frac{1}{q-1} \frac{\int \frac{\frac{\partial}{\partial t} \rho_t^q}{\nu^{q-1}} dx}{F_{q,\nu}(\rho_t)} \\ &= \frac{q}{q-1} \frac{\int \left(\frac{\rho_t}{\nu}\right)^{q-1} \frac{\partial \rho_t}{\partial t} dx}{F_{q,\nu}(\rho_t)} \\ &= \frac{q}{(q-1)F_{q,\nu}(\rho_t)} \int \left(\frac{\rho_t}{\nu}\right)^{q-1} \frac{\partial \rho_t}{\partial t} dx \\ &= \frac{q}{(q-1)F_{q,\nu}(\rho_t)} \int \left(\frac{\rho_t}{\nu}\right)^{q-1} \left(-\nabla \cdot \left(\rho_t \mathbb{E}_{\rho_{0t}}[s(x_0)|x_t = x]\right) + \Delta \rho_t\right) dx \\ &= \frac{q}{(q-1)F_{q,\nu}(\rho_t)} \int -\rho_t \left\langle \nabla \left(\frac{\rho_t}{\nu}\right)^{q-1}, \nabla \log \frac{\rho_t}{\nu} - \mathbb{E}_{\rho_{0t}}[s(x_0)|x_t = x] + \nabla \log \nu \right\rangle dx \quad \text{integrating by parts} \\ &= \frac{q}{(q-1)F_{q,\nu}(\rho_t)} \left(\int -\rho_t \left\langle \nabla \left(\frac{\rho_t}{\nu}\right)^{q-1}, \nabla \log \frac{\rho_t}{\nu} \right\rangle dx \right. \\ &\quad \left. + \int \rho_t \left\langle \nabla \left(\frac{\rho_t}{\nu}\right)^{q-1}, \mathbb{E}_{\rho_{0t}}[s(x_0)|x_t = x] - \nabla \log \nu \right\rangle dx \right) \\ &= \frac{q}{(q-1)F_{q,\nu}(\rho_t)} \left(\underbrace{- \int \nu \left\langle \nabla \left(\frac{\rho_t}{\nu}\right)^{q-1}, \nabla \frac{\rho_t}{\nu} \right\rangle dx}_{A_1} \right. \\ &\quad \left. + \underbrace{\int \rho_t \left\langle \nabla \left(\frac{\rho_t}{\nu}\right)^{q-1}, \mathbb{E}_{\rho_{0t}}[s(x_0)|x_t = x] - \nabla \log \nu \right\rangle dx}_{A_2} \right) \end{aligned}$$

$A_1 = \frac{4(q-1)}{q^2} \mathbb{E}_\nu \left[\left\| \nabla \left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}} \right\|^2 \right]$, because

$$\begin{aligned} \left\langle \nabla \left(\frac{\rho_t}{\nu}\right)^{q-1}, \nabla \frac{\rho_t}{\nu} \right\rangle &= (q-1) \left\langle \left(\frac{\rho_t}{\nu}\right)^{q-2} \nabla \frac{\rho_t}{\nu}, \nabla \frac{\rho_t}{\nu} \right\rangle \\ &= (q-1) \left\langle \left(\frac{\rho_t}{\nu}\right)^{\frac{q-2}{2}} \nabla \frac{\rho_t}{\nu}, \left(\frac{\rho_t}{\nu}\right)^{\frac{q-2}{2}} \nabla \frac{\rho_t}{\nu} \right\rangle \\ &= (q-1) \left\| \frac{2}{q} \nabla \left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}} \right\|^2 \\ &= \frac{4(q-1)}{q^2} \left\| \nabla \left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}} \right\|^2. \end{aligned}$$

Then we compute A_2 . Note that $\nabla\left(\frac{\rho_t}{\nu}\right)^{q-1} = (q-1)\left(\frac{\rho_t}{\nu}\right)^{q-2}\nabla\frac{\rho_t}{\nu} = (q-1)\left(\frac{\rho_t}{\nu}\right)^{\frac{q-2}{2}}\left(\frac{\rho_t}{\nu}\right)^{\frac{q-2}{2}}\nabla\frac{\rho_t}{\nu} = (q-1)\left(\frac{\rho_t}{\nu}\right)^{\frac{q-2}{2}}\frac{2}{q}\nabla\left(\frac{\rho_t}{\nu}\right)^{q/2}$, therefore

$$\begin{aligned} A_2 &= \int \rho_t \langle \nabla\left(\frac{\rho_t}{\nu}\right)^{q-1}, \mathbb{E}_{\rho_{0t}}[s(x_0)|x_t = x] - \nabla \log \nu \rangle dx \\ &= 2\frac{q-1}{q}\mathbb{E}_{\rho_{0t}}\left[\left(\frac{\rho_t}{\nu}\right)^{\frac{q-2}{2}}\langle \nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}, s(x_0) - \nabla \log \nu(x_t) \rangle\right] \\ &= 2\frac{q-1}{q}\mathbb{E}_{\rho_{0t}}\left[\left\langle \left(\frac{\rho_t}{\nu}\right)^{-\frac{1}{2}}\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}, \left(\frac{\rho_t}{\nu}\right)^{\frac{q-1}{2}}(s(x_0) - \nabla \log \nu(x_t)) \right\rangle\right]. \end{aligned}$$

Applying $\langle x, y \rangle \leq \frac{1}{2q}\|x\|^2 + \frac{q}{2}\|y\|^2$ to $x = \left(\frac{\rho_t}{\nu}\right)^{-\frac{1}{2}}\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}$ and $y = \left(\frac{\rho_t}{\nu}\right)^{\frac{q-1}{2}}(s(x_0) - \nabla \log \nu(x_t))$,

$$\begin{aligned} \left\langle \left(\frac{\rho_t}{\nu}\right)^{-\frac{1}{2}}\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}, \left(\frac{\rho_t}{\nu}\right)^{\frac{q-1}{2}}(s(x_0) - \nabla \log \nu(x_t)) \right\rangle &\leq \frac{1}{2q}\left\|\left(\frac{\rho_t}{\nu}\right)^{-\frac{1}{2}}\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}\right\|^2 + \frac{q}{2}\left\|\left(\frac{\rho_t}{\nu}\right)^{\frac{q-1}{2}}(s(x_0) - \nabla \log \nu(x_t))\right\|^2 \\ &= \frac{1}{2q}\frac{\nu}{\rho_t}\left\|\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}\right\|^2 + \frac{q}{2}\left(\frac{\rho_t}{\nu}\right)^{q-1}\|s(x_0) - \nabla \log \nu(x_t)\|^2. \end{aligned}$$

$$\begin{aligned} \implies A_2 &\leq 2\frac{q-1}{q}\left(\frac{1}{2q}\mathbb{E}_{\nu}\left[\left\|\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}\right\|^2\right] + \frac{q}{2}\mathbb{E}_{\rho_{0t}}\left[\left(\frac{\rho_t}{\nu}\right)^{q-1}\|s(x_0) - \nabla \log \nu(x_t)\|^2\right]\right) \\ &= \frac{q-1}{q^2}\mathbb{E}_{\nu}\left[\left\|\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}\right\|^2\right] + (q-1)\mathbb{E}_{\rho_{0t}}\left[\left(\frac{\rho_t}{\nu}\right)^{q-1}\|s(x_0) - \nabla \log \nu(x_t)\|^2\right]. \end{aligned}$$

Therefore

$$\begin{aligned} \frac{\partial}{\partial t}R_{q,\nu}(\rho_t) &= \frac{q}{(q-1)F_{q,\nu}(\rho_t)}(-A_1 + A_2) \\ &\leq \frac{q}{(q-1)F_{q,\nu}(\rho_t)}\left(-\frac{3(q-1)}{q^2}\mathbb{E}_{\nu}\left[\left\|\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}\right\|^2\right] + (q-1)\mathbb{E}_{\rho_{0t}}\left[\left(\frac{\rho_t}{\nu}\right)^{q-1}\|s(x_0) - \nabla \log \nu(x_t)\|^2\right]\right) \\ &= -\frac{1}{F_{q,\nu}(\rho_t)}\left(\frac{3}{q}\mathbb{E}_{\nu}\left[\left\|\nabla\left(\frac{\rho_t}{\nu}\right)^{\frac{q}{2}}\right\|^2\right] - q\mathbb{E}_{\rho_{0t}}\left[\left(\frac{\rho_t}{\nu}\right)^{q-1}\|s(x_0) - \nabla \log \nu(x_t)\|^2\right]\right). \end{aligned}$$

Let $\varphi_t = \frac{\rho_t}{\nu}$ and $\psi_t = \frac{\varphi_t^{q-1}}{\mathbb{E}_{\nu}[\varphi_t^q]} = \frac{\varphi_t^{q-1}}{F_{q,\nu}(\rho_t)}$, then

$$\frac{\partial}{\partial t}R_{q,\nu}(\rho_t) \leq -\frac{3}{4}q\frac{G_{q,\nu}(\rho_t)}{F_{q,\nu}(\rho_t)} + q\mathbb{E}_{\rho_{0t}}[\psi_t(x_t)\|s(x_0) - \nabla \log \nu(x_t)\|^2].$$

□

Proof of Lemma 6. Lemma 5 states

$$\frac{\partial}{\partial t}R_{q,\nu}(\rho_t) \leq -\frac{3}{4}q\frac{G_{q,\nu}(\rho_t)}{F_{q,\nu}(\rho_t)} + q\mathbb{E}_{\rho_{0t}}[\psi_t(x_t)\|s(x_0) - \nabla \log \nu(x_t)\|^2].$$

All we need is to bound $\mathbb{E}_{\rho_{0t}}[\psi_t(x_t)\|s(x_0) - \nabla \log \nu(x_t)\|^2]$.

$$\mathbb{E}_{\rho_{0t}}[\psi_t(x_t)\|s(x_0) - \nabla \log \nu(x_t)\|^2] \leq 2\underbrace{\mathbb{E}_{\rho_{0t}}[\psi_t(x_t)\|s(x_0) - s(x_t)\|^2]}_{A_3} + 2\mathbb{E}_{\rho_t}[\psi_t(x)\|s(x) - \nabla \log \nu(x)\|^2]$$

where

$$\begin{aligned} A_3 &\leq \mathbb{E}_{\rho_{0t}}\left[\psi_t(x_t)(18L_s^2t^2\|s(x_t) - \nabla \log \nu(x_t)\|^2 + 18L_s^2t^2\|\nabla \log \nu(x_t)\|^2 + 6L_s^2t\|z_0\|^2)\right] \quad \text{by Lemma 2} \\ &= 18L_s^2t^2\mathbb{E}_{\rho_t}\left[\psi_t(x)(\|s(x) - \nabla \log \nu(x)\|^2)\right] + 18L_s^2t^2\mathbb{E}_{\rho_t}\left[\psi_t(x)\|\nabla \log \nu(x)\|^2\right] + 6L_s^2tn \\ &= 18L_s^2t^2\mathbb{E}_{\rho_t\psi_t}\left[\|s(x) - \nabla \log \nu(x)\|^2\right] + 18L_s^2t^2\mathbb{E}_{\rho_t\psi_t}\left[\|\nabla \log \nu(x)\|^2\right] + 6L_s^2tn. \end{aligned}$$

So we have

$$\begin{aligned}
\mathbb{E}_{\rho_0 t} [\psi_t(x_t) \|s(x_0) - \nabla \log \nu(x_t)\|^2] &\leq (36L_s^2 t^2 + 2) \mathbb{E}_{\rho_t \psi_t} [\|s(x) - \nabla \log \nu(x)\|^2] \\
&\quad + 36L_s^2 t^2 \mathbb{E}_{\rho_t \psi_t} [\|\nabla \log \nu(x)\|^2] + 12L_s^2 t n \\
&\leq (36L_s^2 t^2 + 2) M_\infty^2 + \underbrace{36L_s^2 t^2 \mathbb{E}_{\rho_t \psi_t} [\|\nabla \log \nu(x)\|^2]}_{A_5} + 12L_s^2 t n.
\end{aligned}$$

By Lemma 16 in Chewi et al. (2022) under the assumption of $\nabla \log \nu$ being L -Lipschitz,

$$\begin{aligned}
A_5 &\leq \mathbb{E}_{\rho_t \psi_t} [\|\nabla \log \frac{\rho_t \psi_t}{\nu}\|^2] + 2nL \\
&= \mathbb{E}_{\rho_t \psi_t} [\|\frac{\nu}{\rho_t \psi_t} \nabla \frac{\rho_t \psi_t}{\nu}\|^2] + 2nL \\
&= \mathbb{E}_{\rho_t \psi_t} [\|\frac{\nu}{\rho_t \psi_t} \frac{1}{F_{q,\nu}(\rho_t)} \nabla \varphi_t^q\|^2] + 2nL \\
&= \int \frac{\nu^2}{\rho_t \psi_t F_{q,\nu}^2(\rho_t)} \|\nabla \varphi_t^q\|^2 dx + 2nL \\
&= \frac{\mathbb{E}_\nu [\frac{1}{\varphi_t^q} \|\nabla \varphi_t^q\|^2]}{F_{q,\nu}(\rho_t)} + 2nL \\
&= \frac{4\mathbb{E}_\nu [\|\nabla \varphi_t^{\frac{q}{2}}\|^2]}{F_{q,\nu}(\rho_t)} + 2nL \quad \text{by } \frac{1}{\varphi_t^q} \|\nabla \varphi_t^q\|^2 = 4\|\nabla \varphi_t^{\frac{q}{2}}\|^2 \\
&= q^2 \frac{G_{q,\nu}(\rho_t)}{F_{q,\nu}(\rho_t)} + 2nL.
\end{aligned}$$

Putting together,

$$\begin{aligned}
\frac{\partial}{\partial t} R_{q,\nu}(\rho_t) &\leq (36L_s^2 t^2 q^2 - \frac{3}{4}q) \frac{G_{q,\nu}(\rho_t)}{F_{q,\nu}(\rho_t)} + (36L_s^2 t^2 + 2) M_\infty^2 q + 72L_s^2 t^2 nLq + 12L_s^2 t nq \\
&\leq -\frac{1}{2}q \frac{G_{q,\nu}(\rho_t)}{F_{q,\nu}(\rho_t)} + (36L_s^2 h^2 + 2) M_\infty^2 q + 72L_s^2 h^2 nLq + 12L_s^2 h nq \quad \text{since } t^2 \leq h^2 \leq \frac{\alpha^2}{144L_s^2 q^2 L^2} \\
&\leq -\frac{\alpha}{q} R_{q,\nu}(\rho_t) + (36L_s^2 h^2 + 2) M_\infty^2 q + 72L_s^2 h^2 nLq + 12L_s^2 h nq
\end{aligned}$$

where the last inequality is from Lemma 5 in Vempala and Wibisono (2019) under the assumption of ν satisfying LSI. Therefore

$$\frac{\partial}{\partial t} e^{\frac{\alpha}{q} t} R_{q,\nu}(\rho_t) \leq e^{\frac{\alpha}{q} t} \left(72L_s^2 nLqh^2 + 12L_s^2 nqh + (36L_s^2 h^2 + 2) M_\infty^2 q \right).$$

Integrating from 0 to h ,

$$\begin{aligned}
e^{\frac{\alpha}{q} h} R_{q,\nu}(\rho_h) &\leq R_{q,\nu}(\rho_0) + \frac{q(e^{\frac{\alpha}{q} h} - 1)}{\alpha} \left(72L_s^2 nLqh^2 + 12L_s^2 nqh + (36L_s^2 h^2 + 2) M_\infty^2 q \right) \\
&\leq R_{q,\nu}(\rho_0) + 2h \left(72L_s^2 nLqh^2 + 12L_s^2 nqh + (36L_s^2 h^2 + 2) M_\infty^2 q \right)
\end{aligned}$$

where the last inequality is because $e^c - 1 \leq 2c$ for $c = \frac{\alpha}{q} h \in (0, 1)$. Rearranging and renaming $\rho_0 \equiv \rho_k, \rho_h \equiv \rho_{k+1}$, we obtain the recursive contraction

$$R_{q,\nu}(\rho_{k+1}) \leq e^{-\frac{\alpha}{q} h} R_{q,\nu}(\rho_0) + 144L_s^2 nLqh^3 + 24L_s^2 nqh^2 + (72L_s^2 h^3 + 4h) M_\infty^2 q.$$

□

Proof of Theorem 2. Applying the recursion in Lemma 6 k times, we have

$$\begin{aligned}
R_{q,\nu}(\rho_k) &\leq e^{-\frac{\alpha}{q}hk} R_{q,\nu}(\rho_0) + \sum_{i=1}^{k-1} e^{-\frac{\alpha}{q}hi} \left(144L_s^2 n L q h^3 + 24L_s^2 n q h^2 + (72L_s^2 h^3 + 4h) M_\infty^2 q \right) \\
&\leq e^{-\frac{\alpha}{q}hk} R_{q,\nu}(\rho_0) + \frac{1}{1 - e^{-\frac{\alpha}{q}h}} \left(144L_s^2 n L q h^3 + 24L_s^2 n q h^2 + (72L_s^2 h^3 + 4h) M_\infty^2 q \right) \\
&\stackrel{(i)}{\leq} e^{-\frac{\alpha}{q}hk} R_{q,\nu}(\rho_0) + \frac{4q}{3\alpha h} \left(144L_s^2 n L q h^3 + 24L_s^2 n q h^2 + (72L_s^2 h^3 + 4h) M_\infty^2 q \right) \\
&\leq e^{-\frac{\alpha}{2}hk} R_{q,\nu}(\rho_0) + \frac{192nLL_s^2 q^2}{\alpha} h^2 + \frac{32nL_s^2 q^2}{\alpha} h + \left(\frac{96L_s^2 h^2 q^2}{\alpha} + \frac{16q^2}{3\alpha} \right) M_\infty^2.
\end{aligned}$$

where (i) is because $1 - e^{-c} \geq \frac{3}{4}c$ for $c \in (0, \frac{1}{4}]$. \square

F Details of Example 1

Proof. Consider ν as target distribution and we estimate $\nabla \log \nu$ by $\nabla \log \mu$. We will show L^p -accuracy goes to 0 as $m \rightarrow \infty$ but the total variation between ν and μ is lower bounded by a positive number.

For convenience, let $\nu_0 = \mathcal{N}(-m, 1)$, $\nu_1 = \mathcal{N}(m, 1)$ and rewrite $\nu = \frac{3}{4}\nu_0 + \frac{1}{4}\nu_1$, $\mu = \frac{1}{2}\nu_0 + \frac{1}{2}\nu_1$. The lower bound of $\text{TV}(\nu, \mu)$ follows from Devroye et al. (2018); Balasubramanian et al. (2022).

Moreover, Proposition 1 in Balasubramanian et al. (2022) calculates

$$\nabla \log \nu - \nabla \log \mu = -m \frac{\nu_0 \nu_1}{2\nu\mu}.$$

Therefore,

$$\begin{aligned}
\mathbb{E}_\nu \left[\|\nabla \log \nu - \nabla \log \mu\|^p \right] &= \frac{m^p}{2^p} \int \nu \frac{\nu_0^p \nu_1^p}{\nu^p \mu^p} \\
&= \frac{m^p}{2^p} \int \frac{\nu_0^p \nu_1^p}{\nu^{p-1} \mu^p} \\
&= \frac{m^p}{2^p} \int \frac{\nu_0^p \nu_1^p}{\left(\frac{3}{4}\nu_0 + \frac{1}{4}\nu_1\right)^{p-1} \left(\frac{1}{2}\nu_0 + \frac{1}{2}\nu_1\right)^p} \\
&= 4^{p-1} m^p \int \frac{\nu_0^p \nu_1^p}{(3\nu_0 + \nu_1)^{p-1} (\nu_0 + \nu_1)^p} \\
&\leq 4^{p-1} m^p \int \frac{\nu_0^p \nu_1^p}{(\nu_0 + \nu_1)^{2p-1}} \quad \text{since } 3\nu_0 + \nu_1 \geq \nu_0 + \nu_1 \\
&\leq 4^{p-1} m^p \left(\int_{\mathbb{R}_-} \frac{\nu_1^p}{\nu_0^{p-1}} + \int_{\mathbb{R}_+} \frac{\nu_0^p}{\nu_1^{p-1}} \right).
\end{aligned}$$

Since

$$\begin{aligned}
\int_{\mathbb{R}_-} \frac{\nu_1^p}{\nu_0^{p-1}} &= \frac{\exp(2p(2p-1)m^2)}{\sqrt{2\pi}} \int_{\mathbb{R}_-} \exp\left(-\frac{1}{2}(x - (2p-1)m)^2\right) dx \\
&= \exp(2p(2p-1)m^2) \mathbb{P}_{N(0,1)}\{Z \leq (2p-1)m\} \\
&= \exp(2p(2p-1)m^2) \mathbb{P}_{N(0,1)}\{Z \geq -(2p-1)m\} \\
&\leq \frac{1}{2} \exp(2p(2p-1)m^2) \exp\left(-\frac{(2p-1)^2 m^2}{2}\right) \quad \text{by Gaussian tail} \\
&= \frac{1}{2} \exp\left(-\frac{m^2}{2}\right).
\end{aligned}$$

and similarly

$$\begin{aligned}\int_{\mathbb{R}_+} \frac{\nu_0^p}{\nu_1^{p-1}} &= \frac{\exp(2p(2p-1)m^2)}{\sqrt{2\pi}} \int_{\mathbb{R}_-} \exp\left(-\frac{1}{2}(x+(2p-1)m)^2\right) dx \\ &= \exp(2p(2p-1)m^2) \mathbb{P}_{N(0,1)}\{Z \geq -(2p-1)m\} \\ &\leq \frac{1}{2} \exp\left(-\frac{m^2}{2}\right).\end{aligned}$$

Therefore,

$$\mathbb{E}_\nu \left[\|\nabla \log \nu - \nabla \log \mu\|^p \right] \leq 4^{p-1} m^p \exp\left(-\frac{m^2}{2}\right) \rightarrow 0 \quad \text{as } m \rightarrow \infty.$$

□