

GOAL

GOAL: Sample from a probability distribution π supported on $\mathcal{X} \subset \mathbb{R}^D$ in a high dimensional setting (i.e., for a large p).

APPLICATIONS: Bayesian inference, generative modeling, etc.

KNOWN: $f \triangleq -\log(\frac{d\pi}{d\mathbf{x}})$, ($f \in C^2(\mathcal{X})$)

(Euclidean) Langevin Monte Carlo

■ (Overdamped) Langevin dynamics:
$$d\mathbf{X}_t = -\nabla f(\mathbf{X}_t)dt + \sqrt{2}dB_t, \quad (\text{LD})$$

where $\{\mathbf{B}_t\}_{t \geq 0}$ is a standard p -dimensional Brownian motion.

■ Euler-Maruyama discretization:
$$\mathbf{X}_{k+1} = \mathbf{X}_k - h_{k+1} \nabla f(\mathbf{X}_k) + \sqrt{2h_{k+1}} \xi_{k+1}, \quad k = 0, 1, 2, \dots \quad (\text{LMC})$$

■ See (Dalalyan and Karagulyan, 2019) for a review

Assumptions: $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle \geq m \|\mathbf{x} - \mathbf{x}'\|_2^2$; (strong convexity)
 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\|_2 \leq M \|\mathbf{x} - \mathbf{x}'\|_2$. (Lipschitz smoothness)

Convergence: Let μ_k be the law of \mathbf{X}_k , $W_2(\cdot, \cdot)$ the Wasserstein 2-distance, and $h_k \equiv h \leq \frac{1}{2(m+M)}$, then

$$W_2(\mu_k, \pi) \leq (1 - m^2 h) W_2(\mu_0, \pi) + 1.65 \left(\frac{M}{m} \right)^{\frac{1}{2}} p^{\frac{1}{2}} h^{\frac{1}{2}}. \quad (1)$$

Iteration Complexity: needs $K_\varepsilon \approx \frac{M^2}{m^2} \log\left(\frac{1}{\varepsilon}\right)$ steps to reach ε -precision.

Examples excluded by the above assumptions

Case 1: Gamma distribution. $f = \sum_{i=1}^p (1 - a_i) \log(x_i) + b_i x_i + C$;

Case 2: Dirichlet distribution. $f = (1 - a_1) \log(x) + (1 - a_2) \log(1 - x) + C$.

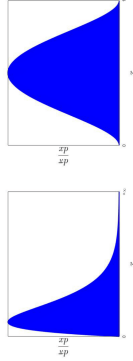


Figure 1: Probability density functions. LEFT: Gamma distribution. RIGHT: Dirichlet distribution.

Relaxation of strong convexity and Lipschitz-smoothness

■ Equivalent assumptions to strong convexity and Lipschitz-smoothness: Let $\phi = \frac{x^2}{2}$,
 $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}'), \nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}') \rangle \geq m \|\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}')\|_2^2$,
 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\|_2 \leq M \|\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}')\|_2. \quad (2)$

■ Relative strong convexity and Lipschitz-smoothness:
 \exists some $C^2(\mathcal{X})$ Legendre-type convex entropy ϕ on \mathcal{X} , such that
 $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}'), \nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}') \rangle \geq m \|\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}')\|_2^2$,
 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\|_2 \leq M \|\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}')\|_2. \quad (3)$

Hessian Riemannian Langevin Monte Carlo algorithm

■ Riemannian Langevin dynamics on Hessian Manifold $(\mathcal{X}, D^2\phi)$ (Roberts and Stramer (2002)):

$$d\mathbf{X}_t = (\theta(\mathbf{X}_t) - [D^2\phi(\mathbf{X}_t)]^{-1} \text{Tr}(\nabla f(\mathbf{X}_t)) dt + \sqrt{2[D^2\phi(\mathbf{X}_t)]^{-1}} dB_t, \quad (4)$$

where $\theta(\mathbf{X}_t) \triangleq -[D^2\phi(\mathbf{X}_t)]^{-1} \text{Tr}([D^3\phi(\mathbf{X}_t)][D^2\phi(\mathbf{X}_t)]^{-1})$.

■ Denoting $\mathbf{Y}_t \triangleq \nabla \phi(\mathbf{X}_t)$, SDE (4) reads

$$d\mathbf{X}_t = -\nabla f \circ \nabla \phi^*(\mathbf{Y}_t) dt + \sqrt{2[D^2\phi^*(\mathbf{Y}_t)]^{-1}} dB_t, \quad (5)$$

here $\phi^*(\mathbf{y}) \triangleq \sup_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \mathbf{y} \rangle - \phi(\mathbf{x})$ is the Legendre-Fenchel conjugate of ϕ .

■ The Euler-Maruyama discretization of SDE (5):

$$\mathbf{Y}_{k+1} \triangleq \mathbf{Y}_k - h_{k+1} \nabla f(\nabla \phi^*(\mathbf{Y}_k)) + \sqrt{2h_{k+1}[D^2\phi^*(\mathbf{Y}_k)]^{-1}} \xi_{k+1}. \quad (6)$$

■ Using $\mathbf{X}_k = \nabla \phi^*(\mathbf{Y}_k)$, we derive the Hessian Riemannian Langevin Monte Carlo algorithm

$$\mathbf{X}_{k+1} \triangleq \nabla \phi^* \left(\nabla \phi(\mathbf{X}_k) - h_{k+1} \nabla f(\mathbf{X}_k) + \sqrt{2h_{k+1}[D^2\phi(\mathbf{X}_k)] \xi_{k+1}} \right). \quad (\text{HRLMC})$$

■ Ignoring the randomness term, HRLMC algorithm reduces to the Mirror Descent algorithm:

$$\mathbf{X}_{k+1} \triangleq \nabla \phi^* \left(\nabla \phi(\mathbf{X}_k) - h_{k+1} \nabla f(\mathbf{X}_k) \right). \quad (\text{Mirror Descent})$$

Other assumptions on ϕ and f

■ Self-concordance-like condition on ϕ :

$$\sqrt{2} \left\| D^2\phi(\mathbf{x})^{\frac{1}{2}} - D^2\phi(\mathbf{x}')^{\frac{1}{2}} \right\|_F \leq \kappa \left\| \nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}') \right\|_2. \quad (7)$$

■ Bound on the commutator of $D^2\phi$ and D^2f :

$$\left\| [D^2\phi(\mathbf{x})]^{-1}, D^2f(\mathbf{x}) \right\|_2 \leq \delta. \quad (8)$$

■ Moment condition on the Hessian of ϕ :

$$R \triangleq \mathbf{E}_{\mathbf{X} \sim \pi} \left\| D^2\phi(\mathbf{X}) \right\|_2 = \int_{\mathcal{X}} \left\| D^2\phi(\mathbf{x}) \right\|_2 e^{-f(\mathbf{x})} d\mathbf{x} < +\infty. \quad (9)$$

■ Interaction of key parameters:

$$\kappa \triangleq \sqrt{\kappa^2 + \frac{\delta(4M + \delta)}{2(m + M)}} < \sqrt{2m}. \quad (10)$$

Note: Our theory covers not only all the strongly convex and Lipschitz-smooth cases, but also new cases not known in the literature.

Wasserstein Distance

■ Let d be the Riemannian distance associated with the squared Hessian metric $[D^2\phi(\mathbf{x})]^2$. Then the natural associated geometric distance on the space of probability distributions on \mathcal{X} is the Wasserstein distance

$$W_{2,\phi}(\mu, \nu) \triangleq \inf_{\mathbf{x} \sim \mu, \mathbf{x}' \sim \nu} \mathbf{E} \left[d^2(\mathbf{x}, \mathbf{x}') \right] = \inf_{\mathbf{x} \sim \mu, \mathbf{x}' \sim \nu} \mathbf{E} \left[\left\| \nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}') \right\|_2^2 \right].$$

Note: When $\phi(\mathbf{x}) = \|\mathbf{x}\|^2/2$, one recovers the standard W_2 distance used in the Euclidean Langevin Monte Carlo (1).

Main Result

Theorem (Constant step-size). Under the assumptions (3),(7)-(10), assume $h_k \equiv h$ is sufficiently small. Then

$$W_{2,\phi}(\mu_k, \pi) \leq \rho^k W_{2,\phi}(\mu_0, \pi) + h^{\frac{1}{2}} (1 - \rho)^{-1} p^{\frac{1}{2}} M R \left(1.65 \sqrt{M} + \kappa / \sqrt{3} \right) + h(1 - \rho)^{-1} p^{\frac{1}{2}} R^{\frac{1}{2}},$$

where the contraction ratio $\rho \triangleq \sqrt{(1 - mh)^2 + h^2} < 1$.

■ **Contraction:** Under vanishing step-sizes, the **HRLMC** algorithm contracts toward a Wasserstein ball centered at the target distribution π with radius

$$r_0 \triangleq \frac{2\kappa p^{\frac{1}{2}} R^{\frac{1}{2}}}{2m - \kappa^2}.$$

Note: when $\phi = \frac{x^2}{2}$, $r_0 = 0$ implies convergence.

■ **Iteration Complexity:** $K_\varepsilon \approx \frac{pMR(\sqrt{M} + \kappa)^2}{(2m - \kappa^2)^2} \frac{1}{\varepsilon^2} \log\left(\frac{1}{\varepsilon}\right)$ steps to reach $(r_0 + \varepsilon)$ -precision.

Numerics

■ Dirichlet distribution $d\pi \propto x^{\alpha-1}(1-x)^{\beta-1}dx$ on 1D Simplex:

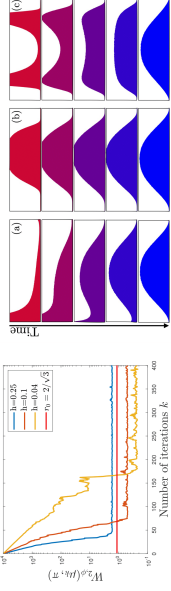


Figure 2: Left: Evolution in time of the sampling error for various constant step-sizes. A horizontal line at $r_0 = \frac{2}{\sqrt{3}}$ materializes the size of the bias term. Right: Visual display of the evolution of the empirical distribution of \mathbf{X}_k for three different initializations.

■ Dirichlet distribution $d\pi \propto x_1^{\alpha_1-1} \dots x_{p-1}^{\alpha_{p-1}-1} (1 - x_1 - \dots - x_{p-1})^{\alpha_p-1} dx_1 \dots dx_{p-1}$ on 2D Simplex:

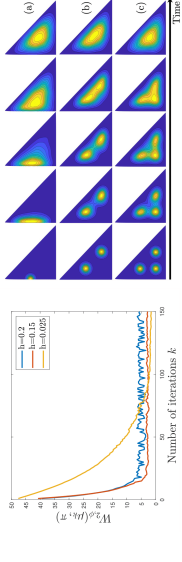


Figure 3: Left: evolution in time of the sampling error for various constant step-sizes. Right: visual display of the evolution of the empirical distribution of \mathbf{X}_k shown as contour plots at different times, for three different initializations.

Conclusion:

■ Contributions:

- (a) First guarantees of HRLMC which show that it contracts into a Wasserstein ball centered at the desired invariant distribution.
- (b) Our method recovers the state-of-the-art non-asymptotic sampling error bounds in Wasserstein distance for the quadratic case.
- (c) Numerics also support our theory.

■ Future work:

- (a) We conjecture that the bias term is inevitable. How to prove it?
- (b) What is a provably good discretization of the Riemannian Langevin dynamics for general manifolds?