

Quantitative convergence analysis of hypocoercive sampling dynamics

Jianfeng Lu (鲁剑锋)

Duke University

jianfeng@math.duke.edu

MATH+IMS Joint Applied Mathematics Colloquium

Chinese University of Hong Kong, September 2020

Joint with

Yu Cao (NYU)

Lihan Wang (Duke)



$$\mathbb{E}_{X \sim \mu} f(X) \approx \frac{1}{N} \sum_{i=1}^N f(X_i)$$

Sampling high dimensional probability distributions is a ubiquitous challenge in many fields:

- computational statistical mechanics;
- machine learning;
- Bayesian statistics;
- high-dimensional PDEs;
- quantum many-body problems;
- ...

A popular approach is Markov chain Monte Carlo: $\{X_i\}$ sampled from a Markov chain $X_{i+1} \sim p(\cdot|X_i)$ with invariant measure μ .

$$\mathbb{E}_{X \sim \mu} f(X) \approx \frac{1}{N} \sum_{i=1}^N f(X_i)$$

Markov chain Monte Carlo: $\{X_i\}$ sampled from a Markov chain $X_{i+1} \sim p(\cdot | X_i)$ with invariant measure μ .

Central limit theorem holds for “nice” Markov chains:

$$\sqrt{N} \left(\frac{1}{N} \sum_{i=1}^N f(X_i) - \mathbb{E}_{\mu} f(X) \right) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

with asymptotic variance

$$\sigma^2 = \text{var}[f(X_i)] + 2 \sum_{k=1}^{\infty} \text{cov}[f(X_i), f(X_{i+k})].$$

Efficiency of the MCMC sampler desires

- Short burn-in period;
- Small asymptotic variance.

Efficiency of the MCMC sampler desires

- Short burn-in period;
- Small asymptotic variance.

This talk: Continuous state space $x \in \mathbb{R}^d$, in particular $d \gg 1$

Common design principle:

Construct a continuous time Markov process and then discretize.

Example: Overdamped Langevin dynamics for $d\mu \propto e^{-U(x)} dx$

$$dx_t = -\nabla U(x_t) dt + \sqrt{2} dW_t.$$

[Rosky, Doll, Friedman 1978]; [Besag 1994]; [Roberts, Tweedie 1996]

Hope for fast convergence to equilibrium of the sampling dynamics.

Overdamped Langevin dynamics for $d\mu \propto e^{-U(x)} dx$

$$dx_t = -\nabla U(x_t) dt + \sqrt{2} dW_t.$$

The Fokker-Planck equation (backward Kolmogorov equation)

$$\partial_t h = -\nabla_x U \cdot \nabla_x h + \Delta_x h, \quad h(0, x) = h_0(x).$$

The convergence of Fokker-Planck equation is well understood, as the generator is self-adjoint and coercive with respect to L^2_μ .

Assumption (Poincaré inequality for μ)

$$\int (h - \int h d\mu)^2 d\mu \leq \frac{1}{m} \int |\nabla_x h|^2 d\mu$$

This implies that the overdamped dynamics has **convergence rate m** .

$$\|h(t, \cdot) - \int h(t, \cdot) d\mu\|_{L^2(\mu)} \leq e^{-mt} \|h(0, \cdot) - \int h(0, \cdot) d\mu\|_{L^2(\mu)}$$

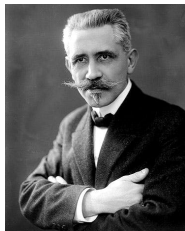
Our motivation is to establish quantitative convergence rate estimate for **hypocoercive** sampling dynamics.

Our first example is the (underdamped) Langevin dynamics

$$dx_t = v_t dt$$

$$dv_t = -\nabla U(x_t) dt - \gamma v_t dt + \sqrt{2\gamma} dW_t$$

Here γ is a friction parameter.



Paul Langevin (1872–1946)

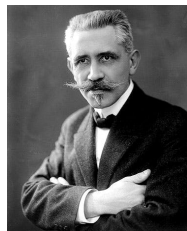
Our motivation is to establish quantitative convergence rate estimate for **hypocoercive** sampling dynamics.

Our first example is the (underdamped) Langevin dynamics

$$dx_t = v_t dt$$

$$dv_t = -\nabla U(x_t) dt - \gamma v_t dt + \sqrt{2\gamma} dW_t$$

Here γ is a friction parameter.



Paul Langevin (1872–1946)

As $\gamma \rightarrow \infty$, and after a time rescaling, we will recover the overdamped Langevin dynamics $dx_t = -\nabla U(x_t) dt + \sqrt{2} dW_t$.

The invariant measure of the Langevin dynamics is given by

$$\rho_\infty(dx, dv) = \frac{1}{Z} e^{-U(x) - \frac{1}{2}|v|^2} dx dv,$$

where Z is the normalizing constant. The marginal distribution is μ .

Langevin dynamics

$$dx_t = v_t dt$$

$$dv_t = -\nabla U(x_t) dt - \gamma v_t dt + \sqrt{2\gamma} dW_t$$

The corresponding backward Kolmogorov equation, known as the kinetic Fokker-Planck equation, is given by

$$\partial_t f = \mathcal{L}f$$

$$f(0, x, v) = f_0(x, v)$$

with the generator given by $\mathcal{L} = \mathcal{L}_{\text{ham}} + \gamma \mathcal{L}_{\text{FD}}$ with

$$\mathcal{L}_{\text{ham}} = v \cdot \nabla_x - \nabla_x U \cdot \nabla_v \quad \text{and} \quad \mathcal{L}_{\text{FD}} = \Delta_v - v \cdot \nabla_v$$

We can verify that $\mathcal{L}^* \rho_\infty = 0$, and thus ρ_∞ the invariant measure.

Recall that the overdamped Langevin dynamics converges with rate m , where m is the Poincaré constant of μ .

Question: Any improvement by the underdamped Langevin dynamics?

Recall that the overdamped Langevin dynamics converges with rate m , where m is the Poincaré constant of μ .

Question: Any improvement by the underdamped Langevin dynamics?

Theorem (Cao-L.-Wang 2019)

For convex U satisfying $|\text{Hess } U| \lesssim (1 + |\nabla U|)$ and superlinear as $|x| \rightarrow \infty$,

$$\|f(t, \cdot) - \int f(t, \cdot) d\rho_\infty\|_{L^2(\rho_\infty)} \leq C_0 \exp(-\lambda t) \|f(0, \cdot) - \int f(0, \cdot) d\rho_\infty\|_{L^2(\rho_\infty)}$$

with explicit estimate of λ as

$$\lambda = \sqrt{m} \log\left(1 + \frac{\gamma\sqrt{m}}{c_0(\sqrt{m} + \gamma)^2}\right)$$

Recall that the overdamped Langevin dynamics converges with rate m , where m is the Poincaré constant of μ .

Question: Any improvement by the underdamped Langevin dynamics?

Theorem (Cao-L.-Wang 2019)

For convex U satisfying $|\text{Hess } U| \lesssim (1 + |\nabla U|)$ and superlinear as $|x| \rightarrow \infty$,

$$\|f(t, \cdot) - \int f(t, \cdot) d\rho_\infty\|_{L^2(\rho_\infty)} \leq C_0 \exp(-\lambda t) \|f(0, \cdot) - \int f(0, \cdot) d\rho_\infty\|_{L^2(\rho_\infty)}$$

with explicit estimate of λ as

$$\begin{aligned} \lambda &= \sqrt{m} \log\left(1 + \frac{\gamma\sqrt{m}}{c_0(\sqrt{m} + \gamma)^2}\right) \\ &= \mathcal{O}(\sqrt{m}) \quad \text{if we take } \gamma = \mathcal{O}(\sqrt{m}). \end{aligned}$$

Results available for more general case; we will not discuss those here.

Exponential convergence with rate \sqrt{m} (setting $\int f \, d\rho_\infty = 0$)

$$\|f(t, \cdot)\|_{L^2(\rho_\infty)} \leq C_0 \exp(-c\sqrt{m}t) \|f(0, \cdot)\|_{L^2(\rho_\infty)}$$

- The $\mathcal{O}(\sqrt{m})$ convergence rate is optimal, as can be seen when U is a Gaussian (so explicit calculation can be done);
- First result in literature for sharp \sqrt{m} convergence rate (acceleration compared with overdamped dynamics with rate m);
- Convergence in L^2 implies convergence of density in χ^2 -divergence, and thus in relative entropy and total variation distance with $\mathcal{O}(\sqrt{m})$ rate.

Our analysis also applies to **piecewise deterministic Markov process**:
Deterministic trajectory between Poisson clocks for random bounces and velocity refreshment

Randomized Hamiltonian Monte Carlo

[Duane, Kennedy, Pendleton, Roweth 1987];
[Bou-Rabee, Sanz-Serna 2017]

$$\mathcal{L} = \underbrace{v \cdot \nabla_x - \nabla_x U \cdot \nabla_v}_{\text{Hamiltonian flow}} + \gamma(\Pi_v - \mathcal{J})$$

Π_v : projection on Gaussian (velocity refreshment)

$$(\Pi_v f)(t, x) := \int f(t, x, v) \kappa(dv)$$

Deterministic Hamiltonian flow in between random (Poisson clock with rate γ) velocity refreshment drawn from Gaussian.

Zigzag sampler (ZZ) [Bierkens, Fearnhead, Roberts 2019]

$$\mathcal{L} = \mathbf{v} \cdot \nabla_{\mathbf{x}} + \sum_{k=1}^d \underbrace{(v_k \partial_{x_k} U)}_{\text{bouncing rate}} + (\mathcal{B}_k - \mathcal{J}) + \gamma(\Pi_{\mathbf{v}} - \mathcal{J})$$

with bouncing operators (flipping the k -th velocity component)

$$\mathcal{B}_k f = f(x, \mathbf{v} - 2v_k \mathbf{e}_k), \quad k = 1, \dots, d$$

Bouncy particle sampler (BPS)

[Peters, de With 2012] [Bouchard-Côté, Vollmer, Doucet 2018]

$$\mathcal{L} = \mathbf{v} \cdot \nabla_{\mathbf{x}} + (\mathbf{v} \cdot \nabla U)_+ (\mathcal{B} - \mathcal{J}) + \gamma(\Pi_{\mathbf{v}} - \mathcal{J})$$

with bouncing operator (flipping wrt hyperplane perpendicular to ∇U)

$$\mathcal{B} f = f\left(x, \mathbf{v} - 2(\mathbf{v} \cdot \nabla U) \frac{\nabla U}{|\nabla U|^2}\right)$$

Promising approaches in the context of stochastic gradient.

Theorem (L.-Wang 2020)

For convex U satisfying $|\text{Hess } U| \lesssim (1 + |\nabla U|)$ and superlinear as $|x| \rightarrow \infty$, all three PDMPs converge exponentially to equilibrium with rates (after an optimal choice of velocity refreshment rate γ)

$$v = \begin{cases} \mathcal{O}(\sqrt{m}), & \text{for RHMC;} \\ \mathcal{O}\left(\frac{\sqrt{m}}{\sqrt{L/m}}\right), & \text{for ZZ;} \\ \mathcal{O}\left(\frac{\sqrt{m}}{\sqrt{d}}\right), & \text{for BPS,} \end{cases}$$

where for the zigzag sampler, we assume in addition that $\nabla^2 U \leq L$.

Results available for more general case; we will not discuss those here.

The rate is optimal for RHMC; for ZZ and BPS, our rate estimate is more quantitative than previous results in [Deligiannidis, Paulin, Bouchard-Côté, Doucet 2018]; [Andrieu, Durmus, Nüsken, Roussel 2018] (which only considered rate dependence in d).

Convergence analysis of kinetic Fokker-Planck equation

$$\partial_t f = \mathcal{L}f = (\mathcal{L}_{\text{ham}} + \gamma \mathcal{L}_{\text{FD}})f; \quad f(0, x, v) = f_0(x, v),$$

where

$$\mathcal{L}_{\text{ham}} = v \cdot \nabla_x - \nabla_x U \cdot \nabla_v \quad \text{and} \quad \mathcal{L}_{\text{FD}} = \Delta_v - v \cdot \nabla_v$$

The operator is not elliptic and only hypo-elliptic [Hörmander 1967] (as the diffusion is degenerate in the x direction).

In particular, we cannot hope for the exponential convergence to follow from a Poincaré (coercivity) estimate for \mathcal{L} . As a result, the constant $C_0 > 1$ is unavoidable in the decay estimate

$$\|f(t, \cdot)\|_{L^2(\rho_\infty)} \leq C_0 \exp(-c\sqrt{mt}) \|f(0, \cdot)\|_{L^2(\rho_\infty)}.$$

Previous results on **quantitative** convergence of Langevin dynamics:

- Convergence in $H_{\rho_\infty}^1$ norm [Villani 2009];
- Convergence in a modified $L_{\rho_\infty}^2$ norm [Dolbeault, Mouhot, Schmeiser 2009; 2015] (also earlier idea from [Herau 2006]).
This was applied to kinetic Fokker-Planck equation by [Roussel, Stoltz 2018], which gives explicit rate estimates, though not sharp
- Very recent result based on resolvent analysis using Schur complement [Bernard, Fathi, Levitt, Stoltz 2020]
- Convergence in Wasserstein distance: using Bakry-Émery framework [Boudoin 2016]; by coupling approaches [Eberle, Guillin, Zimmer 2019; Dalalyan, Riou-Durand 2018]
- Convergence based on Lyapunov function [Mattingly, Stuart, Higham 2002]

Our analysis method was inspired by a recent variational framework [Armstrong, Mourrat 2019], which implicitly used the bracket condition dating back to [Hörmander 1967].

As \mathcal{L} is not coercive, the idea is to resort to augmenting the state space by a time interval $I = (0, T)$ equipped with Lebesgue measure λ . Since in time, the diffusion in v direction will propagate to the x direction.

Our analysis method was inspired by a recent variational framework [Armstrong, Mourrat 2019], which implicitly used the bracket condition dating back to [Hörmander 1967].

As \mathcal{L} is not coercive, the idea is to resort to augmenting the state space by a time interval $I = (0, T)$ equipped with Lebesgue measure λ . Since in time, the diffusion in v direction will propagate to the x direction.

Let κ be the Gaussian measure in velocity ($\rho_\infty(dx dv) = \mu(dx)\kappa(dv)$). The exp. conv. follows from an energy estimate combined with

Theorem (Poincaré inequality in time augmented state space)

$$\|f - (f)_{\lambda \times \mu}\|_{L^2(\lambda \times \mu; L^2_{\kappa})} \lesssim \left(1 + \frac{1}{T\sqrt{m}}\right) \|f - \Pi_v f\|_{L^2(\lambda \times \mu; L^2_{\kappa})} + \left(\frac{1}{\sqrt{m}} + T\right) \|\partial_t f - \mathcal{L}_{ham} f\|_{L^2(\lambda \times \mu; H_{\kappa}^{-1})},$$

where $(f)_{\lambda \times \mu} := \frac{1}{T} \int f(t, x, v) dt d\rho_\infty$.

Proof sketch: Without loss of generality, we assume $(f)_{\lambda \times \mu} = 0$.

By triangular inequality

$$\|f\|_{L^2(\lambda \times \mu; L^2_k)} \leq \|f - \Pi_v f\|_{L^2(\lambda \times \mu; L^2_k)} + \|\Pi_v f\|_{L^2(\lambda \times \mu)}.$$

Proof sketch: Without loss of generality, we assume $(f)_{\lambda \times \mu} = 0$.

By triangular inequality

$$\|f\|_{L^2(\lambda \times \mu; L^2_{\kappa})} \leq \|f - \Pi_v f\|_{L^2(\lambda \times \mu; L^2_{\kappa})} + \|\Pi_v f\|_{L^2(\lambda \times \mu)}.$$

For the underdamped Langevin, using Gaussian Poincaré inequality, we have

$$\|f - \Pi_v f\|_{L^2(\lambda \times \mu; L^2_{\kappa})} \leq \|\nabla_v f\|_{L^2(\lambda \times \mu; L^2_{\kappa})}.$$

The estimate for $\|\Pi_v f\|_{L^2(\lambda \times \mu)}$ is more tricky.

We desire to control it with the help of $\|\partial_t f - \mathcal{L}_{\text{ham}} f\|_{L^2(\lambda \times \mu; H_{\kappa}^{-1})}$, thus, we need to “introduce derivatives”.

The idea is to construct “test functions” $(\phi_0, \phi_1, \dots, \phi_d) \in H_0^1(\lambda \times \mu)$ by solving a divergence equation

$$-\partial_t \phi_0 + \nabla_x^* \cdot \phi = -\partial_t \phi_0 + \sum_{i=1}^d (-\partial_{x_i} \phi_i + \partial_{x_i} U \phi_i) = \Pi_v f$$

with Dirichlet boundary conditions (note that $\int_{I \times \mathbb{R}^d} \Pi_v f \, dt \mu(dx) = 0$).

Lemma

$$\left(\sum_{i=0}^d \|\phi_i\|_{L^2(\lambda \times \mu)}^2 \right)^{1/2} \lesssim \max\{m^{-1/2}, T\} \|\Pi_v f\|_{L^2(\lambda \times \mu)};$$

$$\left(\sum_{i,j=0}^d \|\partial_j \phi_i\|_{L^2(\lambda \times \mu)}^2 \right)^{1/2} \lesssim (1 + m^{-1/2} T^{-1}) \|\Pi_v f\|_{L^2(\lambda \times \mu)}.$$

$$\begin{aligned}
\|\Pi_v f\|_{L^2(\lambda \times \mu)}^2 &= \int_{I \times \mathbb{R}^d} \Pi_v f (-\partial_t \phi_0 + \nabla_x^* \cdot \phi) dt \mu(dx) \\
&= \int_{I \times \mathbb{R}^d \times \mathbb{R}^d} \Pi_v f (-\partial_t \phi_0 + v \cdot \nabla_x \phi_0 \\
&\quad + v \cdot \partial_t \phi - \sum_i v_i v \cdot \partial_{x_i} \phi + \phi \cdot \nabla_x U) dt \rho_\infty(dx dv) \\
&\quad \text{(reintroducing } v \text{ using Gaussianity)}
\end{aligned}$$

After splitting $\Pi_v f$ into f and $\Pi_v f - f$, using integration by parts, we get

$$\begin{aligned}
\|\Pi_v f\|_{L^2(\lambda \times \mu)}^2 &\leq \|\partial_t f - \mathcal{L}_{\text{ham}} f\|_{L^2(\lambda \times \mu; H_k^{-1})} \|\phi_0 - v \cdot \phi\|_{L^2(\lambda \times \mu; H_k^1)} \\
&\quad + \|\!-\partial_t \phi_0 + v \cdot \nabla_x \phi_0 + v \cdot \partial_t \phi - \sum_i v_i v \cdot \partial_{x_i} \phi \\
&\quad + \phi \cdot \nabla_x U\|_{L^2(\lambda \times \mu; L_k^2)} \|f - \Pi_v f\|_{L^2(\lambda \times \mu; L_k^2)}.
\end{aligned}$$

The Poincaré inequality follows from estimate of ϕ and assumption of U .

Quantitative convergence for hypocoercive sampling dynamics based on time-augmented Poincaré inequalities.

- Underdamped Langevin dynamics;
- Randomized Hamiltonian Monte Carlo;
- Zigzag sampler;
- Bouncy particle sampler.

Quantitative convergence for hypocoercive sampling dynamics based on time-augmented Poincaré inequalities.

- Underdamped Langevin dynamics;
- Randomized Hamiltonian Monte Carlo;
- Zigzag sampler;
- Bouncy particle sampler.

We did not discuss in this talk the convergence of the sampling algorithm based on discretization; non-asymptotic analysis of those has been an active research area in machine learning and statistics literature. See e.g.,

- Discretized underdamped Langevin dynamics: [Cheng, Chatterji, Bartlett, Jordan 2018] [Dalalyan, Riou-Durand 2018] [Mou, Ma, Wainwright, Bartlett, Jordan 2019]; [Shen, Lee 2019];
- Discretized Hamiltonian Monte Carlo: [Mangoubi, Vishnoi 2018]; [Lee, Song, Vempala 2018]; [Chen, Vempala 2019]; [Bou-Rabee, Eberle, Zimmer 2020];

Thank you! Any questions?

Email: jianfeng@math.duke.edu

URL: <http://www.math.duke.edu/~jianfeng/>

References:

- with Yu Cao and Lihan Wang, *On explicit L^2 -convergence rate estimate for underdamped Langevin dynamics*, arXiv:1908.04746
- with Lihan Wang, *On explicit L^2 -convergence rate estimate for piecewise deterministic Markov process*, arXiv:2007.14927