

# PDE analysis for sampling dynamics and generative models

**Jianfeng Lu** (鲁剑锋)

Duke University

[jianfeng@math.duke.edu](mailto:jianfeng@math.duke.edu)

IPAM, UCLA, April 2020

Joint with

**Yu Cao** (Duke → Courant Institute)

**Yulong Lu** (Duke → UMass Amherst)

**Lihan Wang** (Duke)



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

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

Popular approaches for sampling:

- Markov chain Monte Carlo, e.g., diffusion based sampling;
- Interacting samplers, e.g., particle filters;
- Generative models, e.g., normalizing flows, generative adversarial networks.

In this talk, we will discuss two recent sampling-related works:

Convergence to equilibrium of Langevin dynamics  
(joint with Yu Cao and Lihan Wang)

Generative models for high dimensional probability distributions  
(joint with Yulong Lu)

Convergence to equilibrium of Langevin dynamics  
(joint with Yu Cao and Lihan Wang)

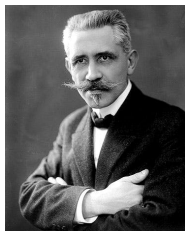
Generative models for high dimensional probability distributions  
(joint with Yulong Lu)

## Langevin dynamics

Consider the (underdamped) Langevin dynamics:

$$\begin{aligned} dx_t &= v_t dt \\ m dv_t &= -\nabla U(x_t) dt - \gamma v_t dt + \sqrt{2\gamma k_B T} dW_t \end{aligned}$$

- $x$ : position;  $v$ : velocity;
- $U$ : potential energy;
- $\gamma$ : friction coefficient;
- $m$ : particle mass;
- $T$ : temperature;
- $k_B$  Boltzmann constant;
- $W_t$ :  $d$ -dimensional standard Brownian motion.



Paul Langevin (1872–1946)

For simplicity, we will set  $m = 1$  and  $k_B T = 1$  (after possibly a rescaling).

Langevin dynamics (with only friction parameter  $\gamma$ )

$$dx_t = v_t dt$$

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

Langevin dynamics (with only friction parameter  $\gamma$ )

$$dx_t = v_t dt$$

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

As  $\gamma \rightarrow 0$ , we get the Hamiltonian dynamics for  $H(v, x) = \frac{1}{2}|v|^2 + U(x)$ :

$$dx_t = v_t dt = \frac{\partial H(x_t, v_t)}{\partial v} dt$$

$$dv_t = -\nabla U(x_t) dt = -\frac{\partial H(x_t, v_t)}{\partial x} dt$$

Langevin dynamics (with only friction parameter  $\gamma$ )

$$dx_t = v_t dt$$

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

As  $\gamma \rightarrow 0$ , we get the Hamiltonian dynamics for  $H(v, x) = \frac{1}{2}|v|^2 + U(x)$ :

$$dx_t = v_t dt = \frac{\partial H(x_t, v_t)}{\partial v} dt$$

$$dv_t = -\nabla U(x_t) dt = -\frac{\partial H(x_t, v_t)}{\partial x} dt$$

As  $\gamma \rightarrow \infty$ , on the  $O(1)$  time scale, it is dominated by the Ornstein-Uhlenbeck process in the velocity variable with "equilibrium"

$$"v_t dt = -\gamma^{-1} \nabla U(x_t) dt + \sqrt{2\gamma^{-1}} dW_t,"$$

after rescaling  $\gamma t \mapsto t$ , we obtain the overdamped Langevin dynamics

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



## 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

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

## 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

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

The invariant measure of the Langevin dynamics is given by

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

where  $Z$  is the normalizing constant (recall that  $\beta = k_B T = 1$ ).

The invariant measure of Langevin dynamics is

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

its marginal on  $x$  is the Boltzmann-Gibbs distribution  $\mu_U \propto e^{-U(x)}$ .

The invariant measure of Langevin dynamics is

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

its marginal on  $x$  is the Boltzmann-Gibbs distribution  $\mu_U \propto e^{-U(x)}$ .

Thus, the Langevin dynamics and its overdamped limit can be used to sample the Boltzmann-Gibbs distribution.

- Smart Monte Carlo [Rosky, Doll, Friedman 1978] (which in fact discussed both underdamped and overdamped Langevin);
- The overdamped version is known in the statistics literature as the Metropolis-adjusted Langevin algorithm (MALA) [Besag 1994; Roberts, Tweedie 1996];
- The underdamped version becomes popular in machine learning in recent years, see e.g., [Cheng, Chatterji, Bartlett, Jordan 2018];
- The underdamped version is also closely related to the Hamiltonian Monte Carlo method [Duane, Kennedy, Pendleton, Roweth 1987];
- Often combined with Metropolis algorithm, though can be used without accept/reject, see e.g., [Neal 1993].

**Question:** How fast Langevin dynamics converges to equilibrium?

Recall the 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

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

Since  $\mathcal{L}^* \rho_\infty = 0$ ,  $\int f(t, \cdot) d\rho_\infty$  is invariant in time. As  $t \rightarrow \infty$ ,  $f$  will converge to the constant  $\int f(0, \cdot) d\rho_\infty$ . **How fast?**

**Question:** How fast Langevin dynamics converges to equilibrium?

Recall the 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

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

Since  $\mathcal{L}^* \rho_\infty = 0$ ,  $\int f(t, \cdot) d\rho_\infty$  is invariant in time. As  $t \rightarrow \infty$ ,  $f$  will converge to the constant  $\int f(0, \cdot) d\rho_\infty$ . **How fast?**

c.f. the Fokker-Planck equation for overdamped dynamics

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

The convergence of Fokker-Planck is well understood, as the generator is self-adjoint and coercive with respect to  $L^2_{\mu_U}$ .

### Assumption (Poincaré inequality for $\mu_U$ )

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

This implies that the overdamped dynamics has **convergence rate  $m$** .  
e.g., when  $U$  is  $m$ -strongly convex.

Question: Any improvement by the underdamped Langevin dynamics?

## Assumption (Poincaré inequality for $\mu_U$ )

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

This implies that the overdamped dynamics has **convergence rate  $m$** .  
e.g., when  $U$  is  $m$ -strongly convex.

Question: Any improvement by the underdamped Langevin dynamics?

gradient descent  $\rightarrow$  accelerated gradient descent

overdamped dynamics  $\rightarrow$  Langevin dynamics

Langevin dynamics can be thought of as a stochastic version of accelerated gradient dynamics.

Thus, if we take the analogy of optimization vs sampling, we would hope that (underdamped) Langevin dynamics gives **convergence rate  $\mathcal{O}(\sqrt{m})$**  for strongly convex  $U$ .



### Assumption (Poincaré inequality for $\mu_U$ )

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

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

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  $\lambda$  as

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

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

Exponential convergence with rate  $\mathcal{O}(\sqrt{m})$  (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  $\chi^2$ -divergence, and thus in relative entropy and total variation distance with  $\mathcal{O}(\sqrt{m})$  rate;
- The constant  $C_0 > 1$  is unavoidable, since the operator  $\mathcal{L}$  is not coercive (more on next slide).

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 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  $\mathcal{L}$ . Hypo-coercivity needs to be established.

## 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 analysis of kinetic Fokker-Planck eqns dates back to [Kolmogorov 1934]. Related previous results on convergence:

- 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 constants, though not sharp.
- Convergence in Wasserstein distance:  
Bakry-Émery framework [Boudoin 2016];  
Coupling approaches [Eberle, Guillin, Zimmer 2019; Cheng, Chatterji, Bartlett, Jordan 2018; Dalalyan, Riou-Durand 2018].

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

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

Our analysis method was inspired by a recent variational framework [Armstrong, Mourrat 2019] for kinetic FP equation on compact domain without potential, 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  $\lambda$ . Since in time, the diffusion in  $v$  direction will propagate to the  $x$  direction.

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

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

Let  $\kappa$  be the Gaussian measure in velocity ( $\rho_\infty(dx dv) = \mu_U(dx)\kappa(dv)$ ).  
The exp. conv. follows from separating into intervals with  $T = 1/\sqrt{m}$ .

**Theorem (Poincaré inequality in time augmented state space)**

$$\begin{aligned} \|f - (f)_{\lambda \times \mu_U}\|_{L^2(\lambda \times \mu_U; L^2_{\bar{\kappa}})} &\lesssim \left(1 + \frac{1}{T\sqrt{m}}\right) \|\nabla_v f\|_{L^2(\lambda \times \mu_U; L^2_{\bar{\kappa}})} \\ &\quad + \left(\frac{1}{\sqrt{m}} + T\right) \|\partial_t f - \mathcal{L}_{\text{ham}} f\|_{L^2(\lambda \times \mu_U; H_{\bar{\kappa}}^{-1})}, \end{aligned}$$

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

Convergence to equilibrium of Langevin dynamics  
(joint with Yu Cao and Lihan Wang)

Generative models for high dimensional probability distributions  
(joint with Yulong Lu)

**Generative models:** Transform simple probability distributions (such as Gaussian) often by neural networks to the desired probability distribution:

$$\pi \approx g_{\#}p_z.$$

**Question?** Given source and target distributions, can one construct a DNN to approximate to push-forward?

This is a distribution version of the question of universal approximation.

cf. [Lee, Ge, Ma, Risteski, Arora 2017]



**Generative models:** Transform simple probability distributions (such as Gaussian) often by neural networks to the desired probability distribution:

$$\pi \approx g_{\#}p_z.$$

**Question?** Given source and target distributions, can one construct a DNN to approximate to push-forward?

This is a distribution version of the question of universal approximation.

cf. [Lee, Ge, Ma, Risteski, Arora 2017]

To make the question more quantitative, we need to specify the metric we use to measure the difference of two distributions.

Integral probability metric (IPM):

$$d_{\mathcal{F}_D}(p, \pi) := \sup_{f \in \mathcal{F}_D} |\mathbb{E}_p f(X) - \mathbb{E}_\pi f(Y)|,$$

where  $f$  can be understood as a discriminator.

Integral probability metric (IPM):

$$d_{\mathcal{F}_D}(p, \pi) := \sup_{f \in \mathcal{F}_D} |\mathbb{E}_p f(X) - \mathbb{E}_\pi f(Y)|.$$

1-Wasserstein distance:  $\mathcal{F}_D$ : 1-Lipschitz functions

$$\mathcal{W}_1(p, \pi) = \inf_{\gamma \in \Gamma(p, \pi)} \int |x - y| \gamma(dx dy).$$

Maximum mean discrepancy (MMD):  $\mathcal{F}_D$ : unit ball of a reproducing kernel Hilbert space (RKHS):

$$\text{MMD}(p, \pi) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} |\mathbb{E}_p f(X) - \mathbb{E}_\pi f(Y)|$$

Kernelized Stein discrepancy (KSD):  $\mathcal{F}_D = \{\mathcal{J}_\pi f \mid f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \leq 1\}$

$$\text{KSD}(p, \pi) = \mathbb{E}_p (\mathcal{J}_\pi f(X))^2$$

where  $\mathcal{J}_\pi$  is the Stein operator (for  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ )

$$\mathcal{J}_\pi f := \nabla \log \pi \cdot f + \nabla \cdot f.$$

## Theorem (L.-Lu 2020)

For source distribution  $p_z$  absolutely continuous wrt Lebesgue measure, there exists a fully connected and feed-forward deep neural network  $u$  with depth  $L = \lceil \log_2 n \rceil$  and width  $N = 2^L$  such that

$$d_{\mathcal{F}_D}((\nabla u)_{\#} p_z, \pi) \leq \varepsilon,$$

where the complexity  $n$  depends on the choice of metric.

- If  $d_{\mathcal{F}_D} = \mathcal{W}_1$  and  $\pi$  has finite 3rd moment,

$$n \lesssim \begin{cases} \varepsilon^{-2}, & d = 1; \\ \log^2(\varepsilon) \varepsilon^{-2}, & d = 2; \\ \varepsilon^{-d}, & d \geq 3. \end{cases}$$

## Theorem (L.-Lu 2020)

For source distribution  $p_z$  absolutely continuous wrt Lebesgue measure, there exists a fully connected and feed-forward deep neural network  $u$  with depth  $L = \lceil \log_2 n \rceil$  and width  $N = 2^L$  such that

$$d_{\mathcal{F}_D}((\nabla u)_{\#} p_z, \pi) \leq \varepsilon,$$

where the complexity  $n$  depends on the choice of metric.

- If  $d_{\mathcal{F}_D} = \text{MMD}$  with kernel  $k$  such that  $\sup_x |k(x, x)| \lesssim 1$ ,

$$n \lesssim \varepsilon^{-2}.$$

- If  $d_{\mathcal{F}_D} = \text{KSD}$  with kernel  $k$  plus some technical assumptions,  $\pi$  is sub-Gaussian and  $\nabla \log \pi(x)$  is Lipschitz,

$$n \lesssim d \varepsilon^{-1}.$$

Remark: It is better to parameterize  $u$  using NN and use  $\nabla u$  to push-forward than directly parameterize the map due to regularity.

**Proof sketch:** The overall strategy is similar to standard approximation theory: Take a sieve for the whole space and construct approximation for elements in the sieve.

- Approximate  $\pi$  by empirical measures  $X_i \sim \pi$ :

$$\pi \approx \pi_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

Wasserstein: [Fourier, Guillin 2015; Weed, Bach 2019; Lei 2020];  
MMD: [Sriperumbudur 2016]. KSD: Our work.

- Push-forward the source dist.  $p_z$  to the empirical distribution  $P_n$ ;  
this is based on the theory of [semi-discrete optimal transport](#).

Monge formulation of optimal transport with quadratic loss

$$\inf_{T: \mathbb{R}^d \rightarrow \mathbb{R}^d} \int \frac{1}{2} |x - T(x)|^2 \mu(dx) \quad \text{s.t. } T_{\#} \mu = \nu.$$

Kantorovich formulation of optimal transport

$$\mathcal{W}_2^2(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \int \frac{1}{2} |x - y|^2 \gamma(dx dy).$$

Monge formulation of optimal transport with quadratic loss

$$T: \mathbb{R}^d \rightarrow \mathbb{R}^d \quad \int \frac{1}{2} |x - T(x)|^2 \mu(dx) \quad \text{s.t. } T_{\#} \mu = \nu.$$

Kantorovich formulation of optimal transport

$$\mathcal{W}_2^2(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \int \frac{1}{2} |x - y|^2 \gamma(dx dy).$$

The Kantorovich dual formulation

$$\begin{aligned} \mathcal{W}_2^2(\mu, \nu) &= \sup_{\substack{\varphi \in L^1(\mu), \psi \in L^1(\nu) \\ \varphi(x) + \psi(y) \leq \frac{1}{2} |x - y|^2}} \int \varphi d\mu + \int \psi d\nu \\ &= \sup_{\psi \in L^1(\nu)} \int \psi^c d\mu + \int \psi d\nu, \end{aligned}$$

where  $\psi^c$  is given by the  $c$ -transform:

$$\psi^c(x) = \inf_{y \in \mathbb{R}^d} \frac{1}{2} |x - y|^2 - \psi(y).$$

Monge formulation of optimal transport with quadratic loss

$$\inf_{T: \mathbb{R}^d \rightarrow \mathbb{R}^d} \int \frac{1}{2} |x - T(x)|^2 \mu(dx) \quad \text{s.t. } T_{\#}\mu = \nu.$$

The Kantorovich dual formulation

$$\mathcal{W}_2^2(\mu, \nu) = \sup_{\psi \in L^1(\nu)} \int \psi^c d\mu + \int \psi d\nu,$$

where  $\psi^c$  is given by the  $c$ -transform:

$$\psi^c(x) = \inf_{y \in \mathbb{R}^d} \frac{1}{2} |x - y|^2 - \psi(y).$$

### Theorem (Brenier 1987)

*Assume that  $\mu$  is absolutely continuous with respect to the Lebesgue measure. Then there exists a convex function  $\bar{\varphi} : \mathbb{R}^d \rightarrow \mathbb{R}$*

$$T(x) = \nabla \bar{\varphi}(x), \quad \mu\text{-a.e. } x.$$

*In fact,  $\bar{\varphi}(x)$  can be chosen as  $\frac{1}{2}|x|^2 - \varphi(x)$ .*



## Semi-discrete optimal transport:

- Source measure  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$  absolutely continuous:  $d\mu = \rho dx$ ;
- Target measure  $\nu = \sum_{j=1}^n \nu_j \delta_{y_j}$ .

In the discrete case, the Kantorovich dual becomes the following functional depending on  $(\psi_1, \dots, \psi_n)$ :

$$\begin{aligned}\mathcal{F}(\psi_1, \dots, \psi_n) &= \int \psi^c(x) \rho(x) dx + \sum_{j=1}^n \psi_j \nu_j \\ &= \int \inf_j \left( \frac{1}{2} |x - y_j|^2 - \psi_j \right) \rho(x) dx + \sum_{j=1}^n \psi_j \nu_j.\end{aligned}$$

## Semi-discrete optimal transport:

- Source measure  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$  absolutely continuous:  $d\mu = \rho dx$ ;
- Target measure  $\nu = \sum_{j=1}^n \nu_j \delta_{y_j}$ .

In the discrete case, the Kantorovich dual becomes the following functional depending on  $(\psi_1, \dots, \psi_n)$ :

$$\begin{aligned}\mathcal{F}(\psi_1, \dots, \psi_n) &= \int \psi^c(x) \rho(x) dx + \sum_{j=1}^n \psi_j \nu_j \\ &= \int \inf_j \left( \frac{1}{2} |x - y_j|^2 - \psi_j \right) \rho(x) dx + \sum_{j=1}^n \psi_j \nu_j.\end{aligned}$$

**Theorem (Gu-Luo-Sun-Yau 2016 for compact domain; L.-Lu 2020)**

*The optimal transport plan  $T$  is given by for some  $m_j \in \mathbb{R}$ .*

$$T(x) = \nabla \bar{\varphi}(x) := \nabla \max_j \{x \cdot y_j + m_j\}.$$

In particular,  $\bar{\varphi}$  is piecewise linear and can be easily approx. by ReLU NN.

- Underdamped Langevin dynamics has faster convergence to equilibrium compared with overdamped dynamics for convex pot.
- Universal approximation without curse-of-dimensionality for generative models for distributions;

Thank you! Any questions?

Email: [jianfeng@math.duke.edu](mailto: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 Yulong Lu, *A universal approximation theorem of deep neural networks for expressing distributions*, arXiv:2004.08867