

Generative Models, Optimal Transport, and PDEs

statistical foundations in the GenAI age

Tengyuan Liang

The University of Chicago

Prelude: motivation

Sensitivity Analysis: denoising diffusions and Fokker-Planck PDE

Interlude: thoughts on generative models

No-Regret Analysis: Sinkhorn algorithm and Monge-Ampère PDE

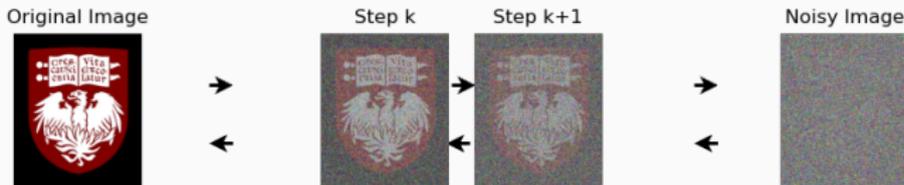
Prelude: motivation

Why do generative models work?

- rethink **learning to sample** from data distributions
- theories for probabilistic generative models

Optimal Transport Perspectives for Generative Models

diffusion-based generative model



Forward chain: **add noise** to **data point**

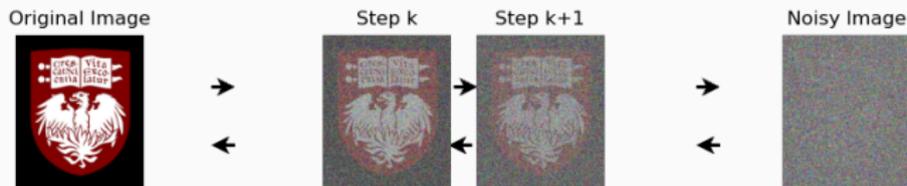
$$X_0 \rightarrow \dots \rightarrow X_k \rightarrow X_{k+1} \rightarrow \dots \rightarrow X_\infty$$

Backward chain: **remove noise** from **data point**

$$X_0 \leftarrow \dots \leftarrow X_k \leftarrow X_{k+1} \leftarrow \dots \leftarrow X_\infty$$

Saremi and Hyvärinen (2019), Sohl-Dickstein, Weiss, Maheswaranathan, and Ganguli (2015), Song and Ermon (2019), and Song, Sohl-Dickstein, et al. (2020)

diffusion-based generative model



Forward chain: Markov chain
diffuse **real data** to **white noise**, **learning** process

Backward chain: time-reversal of a Markov chain
denoise **white noise** to **real data**, **generative/sampling** process

Empirically: state-of-the-art results, DALL-E, Stable Diffusion, ...

Song, Durkan, Murray, and Ermon (2021) and Song and Ermon (2020)

(if infinite computational resources, **perfect** generative model)

paradoxical to a theory researcher

Empirically: state-of-the-art results, DALL-E, Stable Diffusion, ...

Song, Durkan, Murray, and Ermon (2021) and Song and Ermon (2020)

(if infinite computational resources, **perfect** generative model)

However, conceptually problematic in {**probability**, **optimization**, **information**} sense

- **probability**: time-reversal of a Markov chain will get stuck in equilibrium, cannot recover past from future
- **optimization**: if two paths converge to the same point, impossible to recover the initial condition
- **information**: inject noise then denoise, how can one be better off?

~~(if infinite computational resources, **perfect** generative model)~~

resolving the paradox



μ : real data distribution, ν : white noise distribution

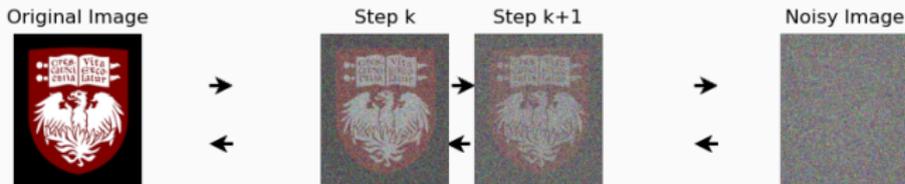
Forward Diffusion $\mu =: \mu_0 \xrightarrow{f_1^\mu} \mu_\eta \rightarrow \dots \xrightarrow{f_k^\mu} \mu_{k\eta} \rightarrow \dots \xrightarrow{f_T^\mu} \mu_{T\eta} \xrightarrow{T \rightarrow \infty} \nu$

Time Reversal $\mu =: \mu_0 \xleftarrow{b_1^\mu} \mu_\eta \leftarrow \dots \xleftarrow{b_k^\mu} \mu_{k\eta} \leftarrow \dots \xleftarrow{b_K^\mu} \mu_{T\eta}$



infinite computation and data: $\nu \xleftarrow{b^\nu} \nu$

resolving the paradox



Time Reversal $\mu =: \mu_0 \xleftarrow{\mathbf{b}_1^\mu} \mu_\eta \xleftarrow{\dots} \mu_{k\eta} \xleftarrow{\mathbf{b}_k^\mu} \dots \xleftarrow{\mathbf{b}_K^\mu} \mu_{K\eta}$



Backward Denoising $\mu \overset{?}{\leftarrow} \bar{\nu}_0 \xleftarrow{\mathbf{b}_1^\mu} \bar{\nu}_\eta \xleftarrow{\dots} \bar{\nu}_{k\eta} \xleftarrow{\mathbf{b}_k^\mu} \dots \xleftarrow{\mathbf{b}_K^\mu} \bar{\nu}_{K\eta} := \nu$

resolving the paradox

Time Reversal $\mu =: \mu_0 \xleftarrow{\mathbf{b}_1^\mu} \mu_\eta \xleftarrow{\dots} \xleftarrow{\mathbf{b}_k^\mu} \mu_{k\eta} \xleftarrow{\dots} \xleftarrow{\mathbf{b}_K^\mu} \mu_{K\eta}$



Backward Denoising $\mu \overset{?}{\leftarrow} \bar{\nu}_0 \xleftarrow{\mathbf{b}_1^\mu} \bar{\nu}_\eta \xleftarrow{\dots} \xleftarrow{\mathbf{b}_k^\mu} \bar{\nu}_{k\eta} \xleftarrow{\dots} \xleftarrow{\mathbf{b}_K^\mu} \bar{\nu}_{K\eta} := \nu$

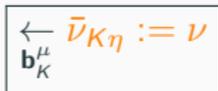
infinite computation and data: denoising maps $\mathbf{b}_k^{\mu'}$'s are exact
denoising maps \mathbf{b}_k^μ carries information of the initial measure μ
(recover the past from the future)

resolving the paradox

Time Reversal $\mu =: \mu_0 \xleftarrow{\mathbf{b}_1^\mu} \mu_\eta \xleftarrow{\dots} \mu_{k\eta} \xleftarrow{\dots} \mu_{K\eta}$

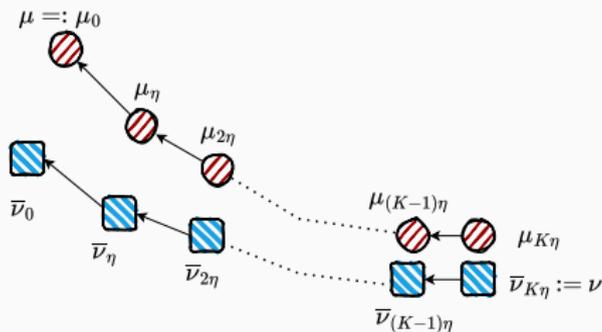


Backward Denoising $\mu \overset{?}{\leftarrow} \bar{\nu}_0 \xleftarrow{\mathbf{b}_1^\mu} \bar{\nu}_\eta \xleftarrow{\dots} \bar{\nu}_{k\eta} \xleftarrow{\dots} \bar{\nu}_{K\eta} := \nu$



Sensitivity analysis

whether the perturbation $\bar{\nu}_{K\eta} \approx \mu_{K\eta}$ will get amplified along the backward denoising chain



what's missing in the current theory

current theory

- sensitivity analysis: on denoising map \mathbf{b}_k^μ , estimated based on finite data, namely $\widehat{\mathbf{b}}_k^\mu \approx \mathbf{b}_k^\mu$
- with infinite data and computation, namely \mathbf{b}_k^μ is known: non-expansion metric $d(\cdot, \cdot)$

$$d(\mu_{k-1}, \bar{\nu}_{k-1}) / d(\mu_k, \bar{\nu}_k) \leq 1$$

f -divergence, data processing inequality

Not fine-grained enough: **the backward denoising chain does not matter**

Song, Durkan, Murray, and Ermon (2021), S. Chen et al. (2022), Lee, Lu, and Tan (2022),
H. Chen, Lee, and Lu (2023), ...

what's missing in the current theory

current theory

with infinite data and computation, namely \mathbf{b}_k^μ is known: non-expansion metric $d(\cdot, \cdot)$

$$d(\mu_{k-1}, \bar{\nu}_{k-1})/d(\mu_k, \bar{\nu}_k) \leq 1$$

Not fine-grained enough: the backward denoising chain does not matter

what's missing in the current theory

current theory

with infinite data and computation, namely \mathbf{b}_k^μ is known: non-expansion metric $d(\cdot, \cdot)$

$$d(\mu_{k-1}, \bar{\nu}_{k-1})/d(\mu_k, \bar{\nu}_k) \leq 1$$

Not fine-grained enough: **the backward denoising chain does not matter**

example: current theory is insufficient

1-dim Gaussian $\mu = \mathcal{N}(0, \sigma^2)$, Wasserstein distance $W(\cdot, \cdot)$

$$W(\mu_{k-1}, \bar{\nu}_{k-1})/W(\mu_k, \bar{\nu}_k) = 1 + \eta \frac{\sigma^2 - 1}{e^{k\eta} + \sigma^2 - 1}$$

$$\begin{cases} > 1, \text{ (expansion)} & \text{if } \sigma^2 > 1 \\ < 1 - \eta, \text{ (strict contraction)} & \text{if } \sigma^2 < \frac{1}{2}, \text{ and } k < \frac{1}{2\eta} \log(2(1 - \sigma^2)) \end{cases}$$

Ideal theory: **the backward denoising chain should matter**

technical challenge: analysis under Wasserstein metric is harder than under KL divergence

S. Chen et al. (2022)

usefulness of the theory: if the message is that backward denoising chain does not matter

Why do diffusion-based generative models work?

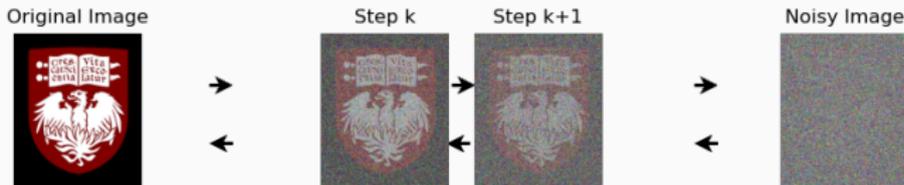
- rethink learning to sample from data distributions
- theories for probabilistic generative models

Optimal Transport Perspectives for Generative Models

a fine-grained theory on the geometric complexity of the
diffuse-then-denoise process

Sensitivity Analysis: denoising diffusions and Fokker-Planck PDE

forward chain: diffusion



On data space, x : Langevin diffusion

$$d\mathbf{X}_t = -\nabla f(\mathbf{X}_t)dt + \sqrt{2\beta^{-1}}d\mathbf{B}_t$$

Lift to distribution of data, $\mu(x)$: Fokker-Planck PDE on $\rho = \frac{d\mu}{dx}$

$$\partial_t \rho_t = \nabla \cdot (\rho_t (\nabla f + \beta^{-1} \nabla \log \rho_t))$$

Jordan-Kinderlehrer-Otto:

$$\mu_{t+\eta} := \arg \min_{\nu \in \mathcal{P}_2^r(X)} \frac{1}{2\eta} W_2^2(\mu_t, \nu) + \mathcal{G}(\nu), \quad \mathcal{G}(\nu) := \int f d\nu + \beta^{-1} \int \log\left(\frac{d\nu}{dx}\right) d\nu$$

$\nabla f + \beta^{-1} \nabla \log \rho$ is the Wasserstein gradient of $\mathcal{G} : \mathcal{P}_2^r(X) \rightarrow \mathbb{R}$

forward chain: diffusion



forward step $\mathbf{X}_t \rightarrow \mathbf{X}_{t+\eta}$ is easy: add Gaussian noise

$$\mathbf{X}_{t+\eta} = (1 - \eta)\mathbf{X}_t + \sqrt{2\beta^{-1}\eta}\mathbf{Z}, \quad \mathbf{Z} \sim \mathcal{N}(0, I)$$

How about backward step? **denoising** $\mathbf{X}_{t+\eta} \rightarrow \mathbf{X}_t$

backward chain: denoising

$$\mathbf{X}_{t+\eta} = (1 - \eta)\mathbf{X}_t + \sqrt{2\beta^{-1}\eta}\mathbf{Z}, \quad \mathbf{Z} \sim \mathcal{N}(0, I)$$

Goal: **denoise** $\mathbf{X}_{t+\eta}$ to predict \mathbf{X}_t

Score function & Eddington-Tweedie formula

Consider $\mathbf{Y} = \mathbf{X} + \sigma\mathbf{Z}$, $\mathbf{Z} \sim \mathcal{N}(0, I_d)$ and \mathbf{X}, \mathbf{Z} are independent. Let $p_{\mathbf{Y}}$ denote the density of \mathbf{Y} . Then

$$\underbrace{\nabla \log p_{\mathbf{Y}}(y)}_{\text{score function}} = \frac{1}{\sigma^2} \left\{ \mathbb{E}[\mathbf{X} | \mathbf{Y} = y] - y \right\}.$$

Score estimation?

$$\min_s \mathbb{E} \left[\left\| \frac{(1 - \eta)\mathbf{X}_t - \mathbf{X}_{t+\eta}}{2\beta^{-1}\eta} - s(\mathbf{X}_{t+\eta}) \right\|^2 \right]$$

Minimizer $s(x)$ is the score function $\nabla \log p_{\mu_{t+\eta}}(x)$

Score function $\nabla \log p_{t+\eta}(x)$ induces a **DenoiseMap** on data level
(based on conditional expectation $\mathbb{E}[\mathbf{X}_t | \mathbf{X}_{t+\eta} = x]$)

However, is the **DenoiseMap** exact on the distribution level?
Namely, $\text{DenoiseMap}(\mu_{t+\eta}) \stackrel{?}{=} \mu_t$

Proposition (L., Koriyama and Dharmakeerthi '24)

$$\mu_{t+\eta} := \arg \min_{\nu \in \mathcal{P}_2^f(X)} \frac{1}{2\eta} W_2^2(\mu_t, \nu) + \mathcal{G}(\nu), \quad \mathcal{G}(\nu) := \int f d\nu + \beta^{-1} \int \log\left(\frac{d\nu}{dx}\right) d\nu.$$

Then any $\eta \in (0, \eta_*)$, the **optimal transport map** $\mathbf{t}_{\mu_{t+\eta}}^{\mu_t}$ is,

$$\mathbf{t}_{\mu_{t+\eta}}^{\mu_t} = \mathbf{i} + \eta(\nabla f + \underbrace{\beta^{-1} \nabla \log p_{\mu_{t+\eta}}}_{\text{score function}}), \text{ for } \mu_{t+\eta}\text{-a.e. } x \in \mathbb{R}^d.$$

$$\mathbf{t}_\nu^\mu := \arg \min_{\mathbf{t} : \mathbf{t}_\# \nu = \mu} \int \|\mathbf{t}(y) - y\|^2 d\nu(y) \quad \text{Monge}$$

$$\begin{aligned} W^2(\mu, \nu) &:= \inf_{\pi \in \Pi(\mu, \nu)} \int \|x - y\|^2 d\pi(x, y) \quad \text{Kantorovich} \\ &= \int \|\mathbf{t}_\nu^\mu - \mathbf{i}\|^2 d\nu \end{aligned}$$

Brenier (1987)

Proposition (L., Koriyama and Dharmakeerthi '24)

$$\mu_{t+\eta} := \arg \min_{\nu \in \mathcal{P}_2^f(X)} \frac{1}{2\eta} W_2^2(\mu_t, \nu) + \mathcal{G}(\nu), \quad \mathcal{G}(\nu) := \int f d\nu + \beta^{-1} \int \log\left(\frac{d\nu}{dx}\right) d\nu .$$

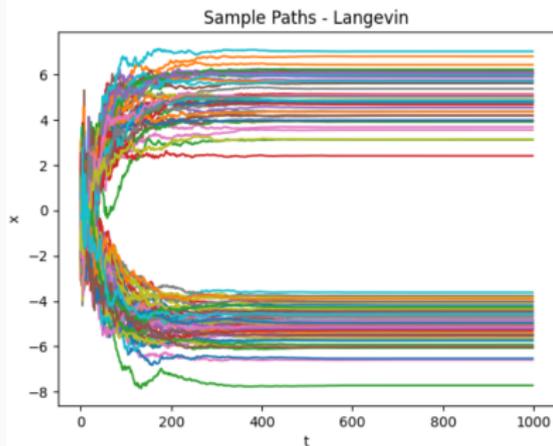
Then any $\eta \in (0, \eta_*)$, the **optimal transport map** $\mathbf{t}_{\mu_{t+\eta}}^{\mu_t}$ is,

$$\mathbf{t}_{\mu_{t+\eta}}^{\mu_t} = \mathbf{i} + \eta(\nabla f + \underbrace{\beta^{-1} \nabla \log p_{\mu_{t+\eta}}}_{\text{score function}}), \text{ for } \mu_{t+\eta}\text{-a.e. } x \in \mathbb{R}^d .$$

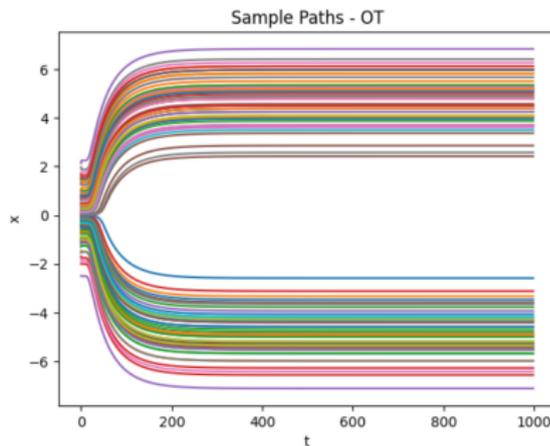
Denoising with score is exact on the **distribution level**, and is **optimal** in Wasserstein sense.

$$X_{(k-1)\eta} = X_{k\eta} + \eta \cdot \left(X_{k\eta} + \beta^{-1} \cdot \widehat{\mathbf{s}}_{k\eta}(X_{k\eta}) \right), \quad \text{ODE}$$

$$X_{(k-1)\eta} = X_{k\eta} + \eta \cdot \left(X_{k\eta} + 2\beta^{-1} \cdot \widehat{\mathbf{s}}_{k\eta}(X_{k\eta}) \right) + \sqrt{2\beta^{-1}\eta} \cdot \mathcal{N}(0, 1), \quad \text{SDE}$$



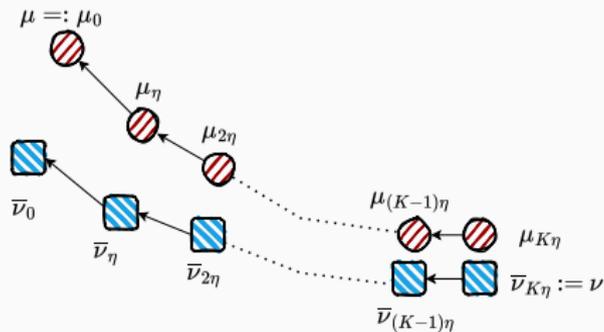
(a) Backward Transport via SDE



(b) Backward Transport via ODE

	statistical	computational	role
score function	conditional mean	optimization step	denoising

sensitivity analysis?



uncertainty in denoising?

Proposition (L., Koriyama and Dharmakeerthi '24)

Consider $\mathbf{Y} = \mathbf{X} + \sigma\mathbf{Z}$, $\mathbf{Z} \sim \mathcal{N}(0, I_d)$ and \mathbf{X}, \mathbf{Z} are independent. Let $p_{\mathbf{Y}}$ denote the density of \mathbf{Y} . Then

$$\underbrace{\nabla^2 \log p_{\mathbf{Y}}(y)}_{\text{curvature function}} = \frac{1}{\sigma^2} \left\{ \underbrace{\text{Cov}\left[\frac{\mathbf{X}}{\sigma} \mid \mathbf{Y} = y\right]}_{\text{denoising uncertainty}} - I_d \right\}.$$

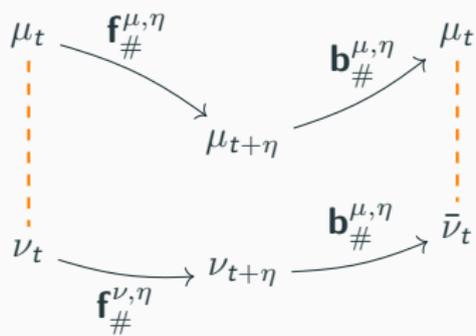
Connects optimization geometry to statistical uncertainty

Turns out crucial in sensitivity analysis

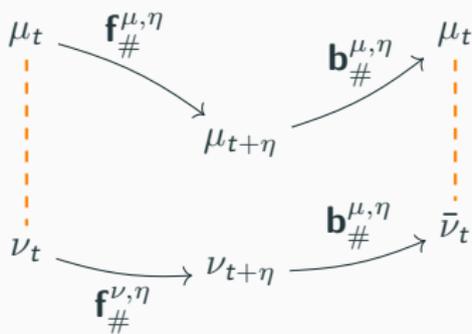
	statistical	computational	role
score function	conditional mean	optimization step	denoising
curvature function	conditional var	optimization guarantee	localization

we discover the crucial role of second-order information in the
diffuse-then-denoise process

warm-up: one-step diffuse-then-denoise



warm-up: one-step **diffuse-then-denoise**



any gain for forward-then-backward step?

Consider $\beta^{-1} = 0$, $\mu_0 = \delta_x$ and $\nu_0 = \delta_y$, two Dirac measures

Forward $W(\mu_0, \nu_0) = \|x - y\|$, Backward $W(\mu_\eta, \bar{\nu}_\eta) = (1 - \eta)\|x - y\|$,

Forward-then-Backward $W(\mu_0, \bar{\nu}_0) = (1 - \eta)^{-1}W(\mu_\eta, \bar{\nu}_\eta) = \|x - y\| = W(\mu_0, \nu_0)$.

No net gain when $\beta^{-1} = 0$.

$$\mu_{t+\eta} := \mathbf{f}_{\#}^{\mu,\eta} \mu_t, \text{ where } \mathbf{f}^{\mu,\eta} = \mathbf{i} - \eta(\nabla f + \beta^{-1} \nabla \log p_{\mu_t})$$

$$\nu_{t+\eta} := \mathbf{f}_{\#}^{\nu,\eta} \nu_t, \text{ where } \mathbf{f}^{\nu,\eta} = \mathbf{i} - \eta(\nabla f + \beta^{-1} \nabla \log p_{\nu_t})$$

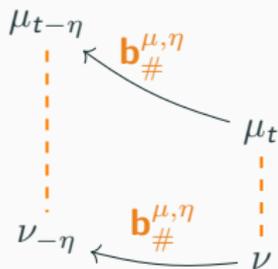
forward contraction

Assume for some $\lambda \in \mathbb{R}_+$, $\nabla^2 f(x) \succeq \lambda \cdot I_d$, $\forall x \in X$.

One-step forward process satisfies

$$\limsup_{\eta \rightarrow 0} \frac{1}{\eta} \frac{W_2^2(\mu_{t+\eta}, \nu_{t+\eta}) - W_2^2(\mu_t, \nu_t)}{W_2^2(\mu_t, \nu_t)} \leq -2\lambda.$$

backward step: sensitivity



$$\mu_{t-\eta} := \mathbf{b}_{\#}^{\mu, \eta} \mu_t,$$

$$\nu_{-\eta} := \mathbf{b}_{\#}^{\mu, \eta} \nu,$$

$$\text{where } \mathbf{b}^{\mu, \eta} = \mathbf{i} + \eta(\nabla f + \beta^{-1} \nabla \log p_{\mu_t})$$

Theorem (L., Koriyama and Dharmakeerthi '24)

Assume for some $\lambda \in \mathbb{R}_+$, $\nabla^2 f(x) \preceq \lambda \cdot I_d$, and for some $\zeta \in \mathbb{R}$

$$\nabla^2 \log p_{\mu_t}(x) \preceq -\zeta \cdot I_d, \quad \forall x \in X.$$

One-step backward map satisfies

$$\sup_{\nu \in \mathcal{P}_2^c(X)} \frac{1}{\eta} \frac{W_2^2(\mathbf{b}_{\#}^{\mu, \eta} \mu_t, \mathbf{b}_{\#}^{\mu, \eta} \nu) - W_2^2(\mu_t, \nu)}{W_2^2(\mu_t, \nu)} \leq 2(\lambda - \beta^{-1} \zeta).$$

The inequality is sharp when μ_t is Gaussian.

diffuse-then-denoise process: one-step

Consider Ornstein–Uhlenbeck process where $f(x) = \|x\|^2/2$.

Forward contraction:

$$\frac{W_2^2(\mu_{t+\eta}, \nu_{t+\eta})}{W_2^2(\mu_t, \nu_t)} \leq 1 - \underline{2\eta} + O(\eta^2).$$

Backward possible expansion:

$$\frac{W_2^2(\mu_t, \bar{\nu}_t)}{W_2^2(\mu_{t+\eta}, \nu_{t+\eta})} \leq 1 + \underline{2\eta} - \boxed{2\eta\beta^{-1}\zeta} + O(\eta^2).$$

diffuse-then-denoise process: one-step

Consider Ornstein–Uhlenbeck process where $f(x) = \|x\|^2/2$.

Forward contraction:

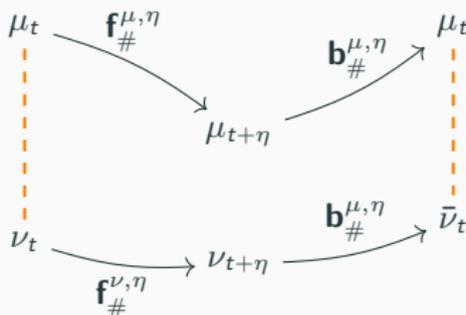
$$\frac{W_2^2(\mu_{t+\eta}, \nu_{t+\eta})}{W_2^2(\mu_t, \nu_t)} \leq 1 - \underline{2\eta} + O(\eta^2).$$

Backward possible expansion:

$$\frac{W_2^2(\mu_t, \bar{\nu}_t)}{W_2^2(\mu_{t+\eta}, \nu_{t+\eta})} \leq 1 + \underline{2\eta} - \boxed{2\eta\beta^{-1}\zeta} + O(\eta^2).$$

Net gain of diffuse-then-denoise:

$$\frac{W_2^2(\mu_t, \mathbf{b}_{\#}^{\mu,\eta} \mathbf{f}_{\#}^{\nu,\eta} \nu_t)}{W_2^2(\mu_t, \nu_t)} \leq 1 - 2\eta\beta^{-1}\zeta + O(\eta^2).$$



curvature $\nabla^2 \log p_{\mu_t}$: **net effect** of diffuse-then-denoise process

curvature at all scales: chaining

Theorem (L., Koriyama and Dharmakeerthi '24), chaining

Define the diffuse-then-noise map

$$\begin{aligned}\mathbf{f}_{[K]}^\nu &:= \mathbf{f}_K^{\nu,\eta} \circ \dots \circ \mathbf{f}_2^{\nu,\eta} \circ \mathbf{f}_1^{\nu,\eta}, \\ \mathbf{b}_{[K]}^\mu &:= \mathbf{b}_1^{\mu,\eta} \circ \mathbf{b}_2^{\mu,\eta} \circ \dots \circ \mathbf{b}_K^{\mu,\eta}.\end{aligned}$$

Assume there exists a sequence of $\zeta_{k\eta} \in \mathbb{R}$ such that

$$\nabla^2 \log p_{\mu_{k\eta}}(x) \preceq -\zeta_{k\eta} \cdot I_d, \quad \forall x \in X, \quad \forall k = 1, 2, \dots, K.$$

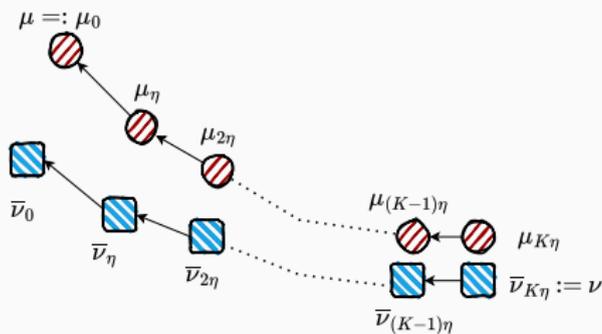
Then the diffuse-then-noise process on ν with fixed K steps satisfies

$$\frac{W_2^2(\mu, (\mathbf{b}_{[K]}^\mu \circ \mathbf{f}_{[K]}^\nu)_\# \nu)}{W_2^2(\mu, \nu)} \leq \exp\left(-2\eta\beta^{-1} \sum_{k=1}^K \zeta_{k\eta} + O(\eta^2)\right).$$

our theory: new insights

This paper: effectiveness of diffuse-then-denoise process

- Net gain of diffuse-then-denoise process: $\beta^{-1} \neq 0$
- Curvature at all scales matters
- If log-concave, $\zeta_{k\eta} > 0$ across all time scales k



This paper: effectiveness of diffuse-then-denoise process

- Net gain of diffuse-then-denoise process: $\beta^{-1} \neq 0$
- Curvature at all scales matters
- If log-concave, $\zeta_{k\eta} > 0$ across all time scales k
- Beyond log-concavity: **worst-case curvature** \rightarrow **average notion of curvature**

multi-scale complexity: beyond log-concavity

Interlude: thoughts on generative models

discussion and future work

	hard to sample	easy to sample
hard to learn	avoid	generative adversarial networks
easy to learn	diffusion-based models	?

Static analysis: Empirical Processes and Optimal Transport

- Generative adversarial networks: one-step push-forward, complexity tradeoff between generator and discriminator

Liang (2021) JMLR

- Optimal transport and generative models: (reversible)-Gromov-Wasserstein, forward and backward maps, isomorphism, cycle-consistency

Hur, Guo, and Liang (2024) SIMODS

Dynamic analysis: Partial Differential Equations and Optimal Transport

- **Fokker-Planck PDE**: probabilistic diffusion models, forward and backward chains, curvature matters

Liang, Dharmakeerthi, and Koriyama (2024)

- **Monge-Ampère PDE**: leverage curvature to design a new chain/flow, regret analysis via new evolution variational inequality

Deb and Liang (2025)

Guo, Hur, Liang, and Ryan (2022) COLT

No-Regret Analysis: Sinkhorn algorithm and Monge-Ampère PDE

- Given target $\frac{d\mu}{dx} = e^{-f}$, and reference $\frac{d\nu}{dx} = e^{-g}$
Find a **convex function** $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$

$$(\nabla\psi)_{\#}\nu = \mu$$

- (Static) Monge-Ampère equation:

$$-f(\nabla\psi(y)) + g(y) + \log \det(\nabla^2\psi(y)) = 0 .$$

- Parabolic Monge-Ampère equation:

$$\frac{\partial\psi_t}{\partial t}(y) = -f(\nabla\psi_t(y)) + g(y) + \log \det(\nabla^2\psi_t(y)) .$$

continuous time limit of Sinkhorn algorithm

Define the density at step k as $\rho_k := (\nabla\psi_k)\#e^{-g}$ and consider

$$\frac{\nabla\psi_{k+1} - \nabla\psi_k}{\eta_k} = -\xi_k, \quad \text{where } \xi_k := \nabla(\log(\rho_k/e^{-f}) \circ \nabla\psi_k).$$

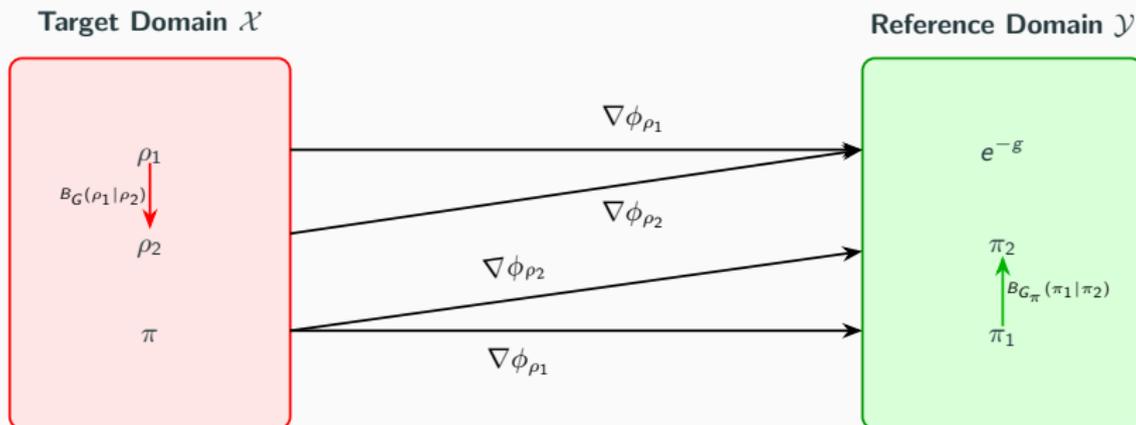
Theorem (Deb and L. '25), evolution variational inequality

Assume that $m_k := \inf_y \nabla^2\psi_k(y) > 0$. Given any probability density π , let π_k denote pdf of $(\nabla\psi_k^*)\#\pi$. Then

$$\begin{aligned} & KL(\rho_k|e^{-f}) - KL(\pi|e^{-f}) \\ &= \boxed{\frac{B_G(\pi|\rho_k) - B_G(\pi|\rho_{k+1})}{\eta_k} + \frac{B_{G_\pi}(\pi_k|\pi_{k+1})}{\eta_k}} - KL(\pi|\rho_k) \\ &\leq \frac{B_G(\pi|\rho_k) - B_G(\pi|\rho_{k+1})}{\eta_k} + \frac{1}{2} \frac{\eta_k}{m_{k+1}} \int \|\xi_k\|^2 d\pi_k - KL(\pi|\rho_k). \end{aligned}$$

Theorem (Deb and L. '25), evolution variational inequality

$$\begin{aligned}
 & KL(\rho_k | e^{-f}) - KL(\pi | e^{-f}) \\
 &= \boxed{\frac{B_G(\pi | \rho_k) - B_G(\pi | \rho_{k+1})}{\eta_k} + \frac{B_{G_\pi}(\pi_k | \pi_{k+1})}{\eta_k}} - KL(\pi | \rho_k) \\
 &\leq \frac{B_G(\pi | \rho_k) - B_G(\pi | \rho_{k+1})}{\eta_k} + \frac{1}{2} \frac{\eta_k}{m_{k+1}} \int \|\xi_k\|^2 d\pi_k - KL(\pi | \rho_k).
 \end{aligned}$$



Theorem (Deb and L. '25), convergence

Average-iterate convergence:

- Constant time-invariant stepsize $\eta_k \equiv \eta$,

$$KL(\bar{\rho}_T | e^{-f}) = O\left(\frac{1}{T}\right)$$

Theorem (Deb and L. '25), convergence

Last-iterate convergence:

- Time-invariant stepsize $\eta_k \equiv \frac{1}{\sqrt{T}}$, assume $\inf_y \nabla^2 \psi_k(y) \geq m \cdot I_d$,

$$KL(\rho_T | e^{-f}) = O\left(\frac{1}{\sqrt{T}}\right)$$

- Adaptive stepsize $\eta_k = \frac{\lambda}{k+1}$, assume in addition $\sup_y \nabla^2 \psi_k(y) \leq M \cdot I_d$, and reference $g(\cdot)$ λ -strongly convex,

$$KL(\rho_T | e^{-f}) = O\left(\frac{\log(T)}{T}\right)$$

- Time-invariant stepsize $\eta_k \equiv \frac{\log(T)}{T}$,

$$B_G(e^{-f} | \rho_T) = O\left(\frac{\log(T)}{T}\right)$$

Monge-Ampère neural-PDE sampler

Algorithm 1: Monge-Ampère Neural-PDE via Logistic Regression

Input : T , total number of steps; $\{\eta_k\}_{k \leq T}$, the step-sizes; initialize a neural network function $\psi_0 : \mathbb{R}^d \rightarrow \mathbb{R}$, and set $k = 0$.

Output: A neural network function $\psi_T : \mathbb{R}^d \rightarrow \mathbb{R}$, and samples from $(\nabla\psi_T)\#e^{-g}$.

while $k < T$ **do**

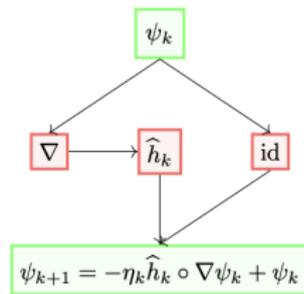
Sampling step: Given i.i.d. data $X_1, \dots, X_n \sim e^{-f}$, sample $\tilde{X}_1, \dots, \tilde{X}_n \sim \rho_k = (\nabla\psi_k)\#e^{-g}$, augment to $\{(X_i, L_i = 0)\}_{i=1}^n \cup \{(\tilde{X}_i, L_i = 1)\}_{i=1}^n$, and form its empirical distribution as $(X, L) \sim \hat{\gamma}_k$;

Learning step: Estimate a neural network discriminator function \hat{h}_k as in (5.1) ;

Neural Network update: Define a new neural network with the residual architecture shown on the right

$\psi_{k+1} = -\eta_k \hat{h}_k \circ \nabla\psi_k + \psi_k$;
Set $k \leftarrow k + 1$;

end



Monge-Ampère neural-PDE sampler

Algorithm 2: Monge-Ampère Neural-PDE via Score Matching

Input : T , total number of steps; $\{\eta_k\}_{k \leq T}$, the step-sizes;
initialize a vector-valued neural network $\mathbf{n}_0 : \mathbb{R}^d \rightarrow \mathbb{R}^d$,
and set $k = 0$.

Output: A vector-valued neural network function $\mathbf{n}_T : \mathbb{R}^d \rightarrow \mathbb{R}^d$,
and samples from $(\mathbf{n}_T) \# e^{-g}$.

while $k < T$ **do**

Sampling step: Obtain i.i.d. samples from $\rho_k = (\mathbf{n}_k) \# e^{-g}$,
 form the empirical measure $\hat{\rho}_k$;

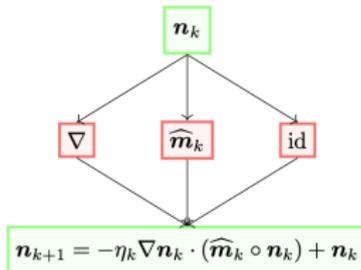
Learning step: Estimate a neural network score $\hat{\sigma}_k$ as in
 (5.2), and define the corresponding $\hat{\mathbf{m}}_k$ as in (5.3) ;

Neural Network update: Define a new neural network with
 the residual architecture shown on the right

$\mathbf{n}_{k+1} = -\eta_k \nabla \mathbf{n}_k \cdot (\hat{\mathbf{m}}_k \circ \mathbf{n}_k) + \mathbf{n}_k$, here $\nabla \mathbf{n}_k : \mathbb{R}^d \rightarrow \mathbb{S}_+^{d \times d}$, and
 $\hat{\mathbf{m}}_k \circ \mathbf{n}_k : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and ‘ \cdot ’ denotes the matrix-vector product ;

 Set $k \leftarrow k + 1$;

end



Thank you!



Tengyuan Liang, Kulunu Dharmakeerthi, and Takuya Koriyama (2024). [Denoising Diffusions with Optimal Transport: Localization, Curvature, and Multi-Scale Complexity.](#)
[arXiv:2411.01629](#)



Nabarun Deb, and Tengyuan Liang (2025). [No-Regret Generative Modeling via Parabolic Monge-Ampère PDE.](#)
[arXiv:2504.09279](#)

multi-scale complexity: beyond log-concavity

SNR, localization, and tail

Definition (L., Koriyama and Dharmakeerthi '24), integrated tail

Given a measure $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, define for any SNR $r \in \mathbb{R}_{\geq 0}$

$$\mathbf{Y}_r := r\mathbf{X} + \mathbf{Z}, (\mathbf{X}, \mathbf{Z}) \sim \mu \otimes \mathcal{N}(0, I_d)$$

SNR, localization, and tail

Definition (**L**., Koriyama and Dharmakeerthi '24), integrated tail

Given a measure $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, define for any SNR $r \in \mathbb{R}_{\geq 0}$

$$\mathbf{Y}_r := r\mathbf{X} + \mathbf{Z}, \quad (\mathbf{X}, \mathbf{Z}) \sim \mu \otimes \mathcal{N}(\mathbf{0}, I_d)$$

Define the localization function, and the associated random variable

$$L_r(y) = \| \text{Cov}[r\mathbf{X} | \mathbf{Y}_r = y] \|_{op}, \quad \text{and } \mathbf{L}_r = \| \text{Cov}[r\mathbf{X} | \mathbf{Y}_r] \|_{op}.$$

and denote the survival function of \mathbf{L}_r as $s_r(\cdot) : \mathbb{R}_+ \rightarrow [0, 1]$

$$s_r(u) := \mathbb{P}(\mathbf{L}_r > u).$$

SNR, localization, and tail

Definition (L., Koriyama and Dharmakeerthi '24), integrated tail

Given a measure $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, define for any SNR $r \in \mathbb{R}_{\geq 0}$

$$\mathbf{Y}_r := r\mathbf{X} + \mathbf{Z}, \quad (\mathbf{X}, \mathbf{Z}) \sim \mu \otimes \mathcal{N}(0, I_d)$$

Define the localization function, and the associated random variable

$$L_r(y) = \|\text{Cov}[r\mathbf{X} | \mathbf{Y}_r = y]\|_{op}, \quad \text{and } \mathbf{L}_r = \|\text{Cov}[r\mathbf{X} | \mathbf{Y}_r]\|_{op}.$$

and denote the survival function of \mathbf{L}_r as $s_r(\cdot) : \mathbb{R}_+ \rightarrow [0, 1]$

$$s_r(u) := \mathbb{P}(\mathbf{L}_r > u).$$

integrated tail:
$$h_\mu(\delta, r) := \int_{1-\delta}^{\infty} s_r(u) du.$$

multi-scale SNR

Ornstein-Uhlenbeck Process with initialization $X_0 \sim \mu$, and potential $f(x) = \|x\|^2/2$

$$d\mathbf{X}_t = -\mathbf{X}_t dt + \sqrt{2\beta^{-1}} d\mathbf{B}_t .$$

Then the distribution of $\mathbf{X}_t, \forall t \in \mathbb{R}_+$

$$\mathbf{X}_t \stackrel{\mathcal{L}}{\sim} e^{-t}\mathbf{X}_0 + \sqrt{\beta^{-1}(1 - e^{-2t})}\mathbf{Z}, \mathbf{Z} \sim \mathcal{N}(0, I_d) .$$

multi-scale SNR

Ornstein-Uhlenbeck Process with initialization $X_0 \sim \mu$, and potential $f(x) = \|x\|^2/2$

$$d\mathbf{X}_t = -\mathbf{X}_t dt + \sqrt{2\beta^{-1}} d\mathbf{B}_t.$$

Then the distribution of $\mathbf{X}_t, \forall t \in \mathbb{R}_+$

$$\mathbf{X}_t \stackrel{\mathcal{L}}{\sim} e^{-t}\mathbf{X}_0 + \sqrt{\beta^{-1}(1 - e^{-2t})}\mathbf{Z}, \quad \mathbf{Z} \sim \mathcal{N}(0, I_d).$$

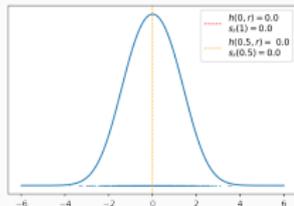
multi-scale SNR

$$\underbrace{\nabla^2 \log p_{\mu_t}(x)}_{\text{curvature function}} = -\frac{1}{s^2(t)} \{I_d - \text{Cov}[r(t)\mathbf{X}_0 | \mathbf{X}_t = x]\}$$

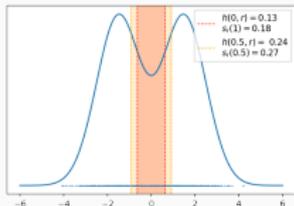
$$r(t) := \frac{e^{-t}}{\sqrt{\beta^{-1}(1 - e^{-2t})}}, \quad s(t) := \frac{1}{\sqrt{\beta^{-1}(1 - e^{-2t})}}$$

simplest non-log-concave example: intuition

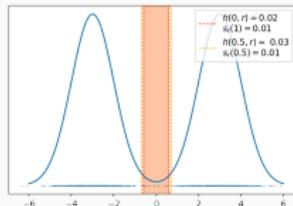
$$\mu = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_{1}$$



(a) $r = 0.71$



(b) $r = 1.50$

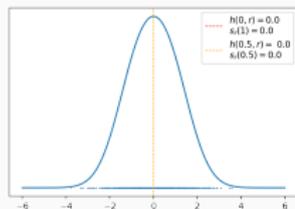


(c) $r = 3.00$

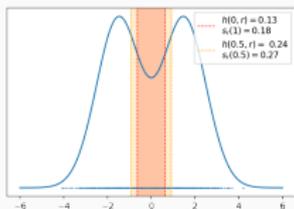
We plot the density $p_{\mathbf{Y}_r}(\cdot)$, for three SNR r 's. Red shaded area corresponds to non-log-concave region with $\nabla^2 \log p_{\mathbf{Y}_r}(\cdot) > -\delta$ with $\delta = 0$, and Orange shaded area corresponds to $\delta = 0.5$.

simplest non-log-concave example: intuition

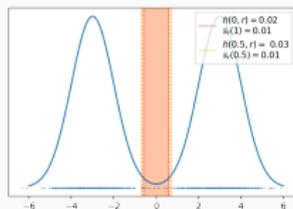
$$\mu = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_{1}$$



(a) $r = 0.71$



(b) $r = 1.50$



(c) $r = 3.00$

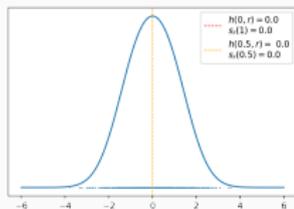
Localization and Curvature: $L_r(y) \leq 1 - \delta$ iff $\nabla^2 \log p_{\mathbf{Y}_r}(y) \preceq -\delta \cdot I_d$, namely, $p_{\mathbf{Y}_r}(y)$ strongly log-concave at y .

These y 's are good locations with strong negative curvature: accurate denoising/localization from y .

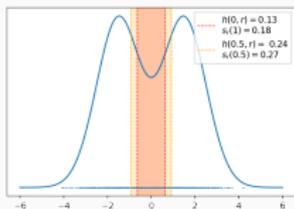
Locations where diffuse-then-denoise is beneficial.

simplest non-log-concave example: intuition

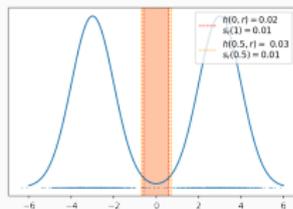
$$\mu = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_{1}$$



(a) $r = 0.71$



(b) $r = 1.50$



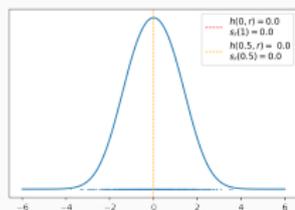
(c) $r = 3.00$

Survival Function: $s_r(1 - \delta) := \mathbb{P}(L_r(\mathbf{Y}_r) > 1 - \delta)$ tells us the probability of bad locations with possibly non-log-concavity where the backward denoising is hard.

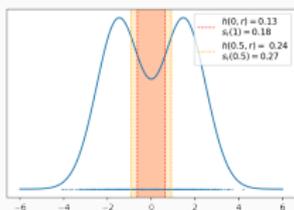
Mass that may induce a large expansion in the diffuse-then-denoise process.

simplest non-log-concave example: intuition

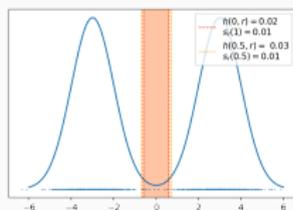
$$\mu = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_{1}$$



(a) $r = 0.71$



(b) $r = 1.50$



(c) $r = 3.00$

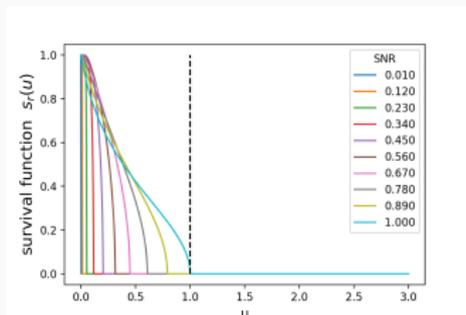
Integrated Tail: slow growth in the integrated tail function $\delta \rightarrow h_{\mu}(\delta, r)$ implies that one can take an effectively large δ

- bad locations with positive curvatures \Rightarrow a negligible expansion effect
- good locations with negative curvature \Rightarrow induce a contraction effect to offset

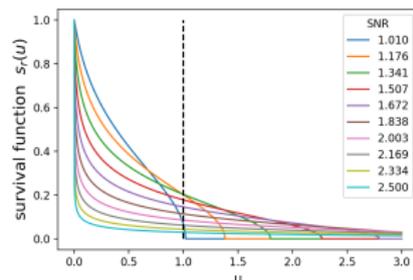
This complexity quantifies an overall notion of curvature.

simplest non-log-concave example: intuition

$$\mu = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_{-1}$$



(a) $s_r(u)$ Low SNR



(b) $s_r(u)$ High SNR

Question: what's the hardest SNR?

multi-scale complexity: theory for non-log-concave

Theorem (L., Koriyama and Dharmakeerthi '24), non-log-concave

For any $\delta \in [0, 1]$, as $\eta \rightarrow 0$,

$$\limsup_{\nu \in \mathcal{M}(\mu_t, M): \nu \xrightarrow{W_2} \mu_t} \frac{1}{\eta} \frac{W_2^2(\mathbf{b}_{\#}^{\mu, \eta} \mu_t, \mathbf{b}_{\#}^{\mu, \eta} \nu) - W_2^2(\mu_t, \nu)}{W_2^2(\mu_t, \nu)} \leq 2 - \frac{2}{1 - e^{-2t}} \left[\underbrace{\delta}_{\text{contraction}} - M \cdot \underbrace{h_{\mu}(\delta, r(t))}_{\text{integrated tail}} \right].$$

multi-scale complexity: theory for non-log-concave

Theorem (L., Koriyama and Dharmakeerthi '24), non-log-concave

For any $\delta \in [0, 1]$, as $\eta \rightarrow 0$,

$$\begin{aligned} \limsup_{\nu \in \mathcal{M}(\mu_t, M) : \nu \xrightarrow{W_2} \mu_t} \frac{1}{\eta} \frac{W_2^2(\mathbf{b}_{\#}^{\mu, \eta} \mu_t, \mathbf{b}_{\#}^{\mu, \eta} \nu) - W_2^2(\mu_t, \nu)}{W_2^2(\mu_t, \nu)} \\ \leq 2 - \frac{2}{1 - e^{-2t}} \left[\underbrace{\delta}_{\text{contraction}} - M \cdot \underbrace{h_{\mu}(\delta, r(t))}_{\text{integrated tail}} \right]. \end{aligned}$$

$$\mathcal{M}(\mu, M) := \left\{ \nu \in \mathcal{P}_2^r(X) : \frac{\sup_{x \in \text{Dom}(\mu)} \|(\mathbf{t}_{\mu}^{\nu} - \mathbf{i})(x)\|^2}{\int \|(\mathbf{t}_{\mu}^{\nu} - \mathbf{i})(x)\|^2 d\mu} \leq M \right\}.$$

Here M controls the richness of perturbations around μ .

multi-scale complexity: theory for non-log-concave

Theorem (L., Koriyama and Dharmakeerthi '24), non-log-concave

For any $\delta \in [0, 1]$, as $\eta \rightarrow 0$,

$$\limsup_{\nu \in \mathcal{M}(\mu_t, M): \nu \xrightarrow{W_2} \mu_t} \frac{1}{\eta} \frac{W_2^2(\mathbf{b}_{\#}^{\mu, \eta} \mu_t, \mathbf{b}_{\#}^{\mu, \eta} \nu) - W_2^2(\mu_t, \nu)}{W_2^2(\mu_t, \nu)} \leq 2 - \frac{2}{1 - e^{-2t}} \left[\underbrace{\delta}_{\text{contraction}} - M \cdot \underbrace{h_{\mu}(\delta, r(t))}_{\text{integrated tail}} \right].$$

probe the multi-scale complexity: empirics

empirics

probing the multi-scale complexity

$$s_r(u) = \mathbb{P}(\mathbf{L}_r > u),$$

$$m^*(r) = \min_{\delta \in [0,1]} \frac{h_\mu(\delta, r)}{\delta},$$

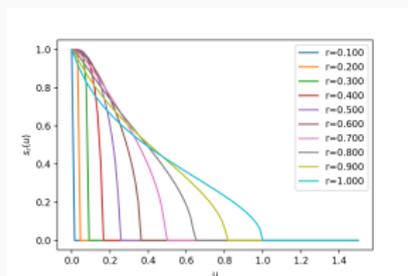
$$\zeta_M^*(t) = \sup_{\delta \in [0,1]} \frac{1}{1-e^{-2t}} [\delta - M \cdot h_\mu(\delta, r(t))].$$

Actionable theory: given a target distribution, empirical or with explicit form, we can visualize the functions above.

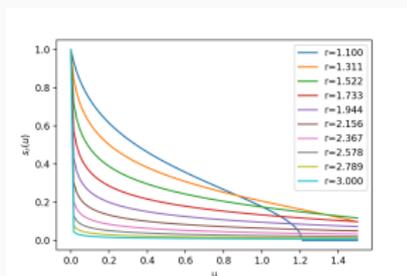
- For any empirical measure μ_0 , a analytic expression for the empirical version of $\mathbf{L}_r(y)$.
- Simulate \mathbf{Y}_r by sampling $r\mathbf{X} + \mathcal{N}(0, 1)$, for $\mathbf{X} \sim \mu_0$.
- Obtain the empirical survival function, $s_r(u)$, $m^*(r)$, and $\zeta_M^*(t)$.

non-log-concave example 1: two-point mass

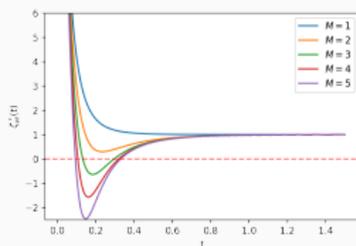
$$X_0 \sim \frac{1}{2}\delta_{-\mu} + \frac{1}{2}\delta_{\mu} \text{ for some } \mu > 0$$



(a)



(b)

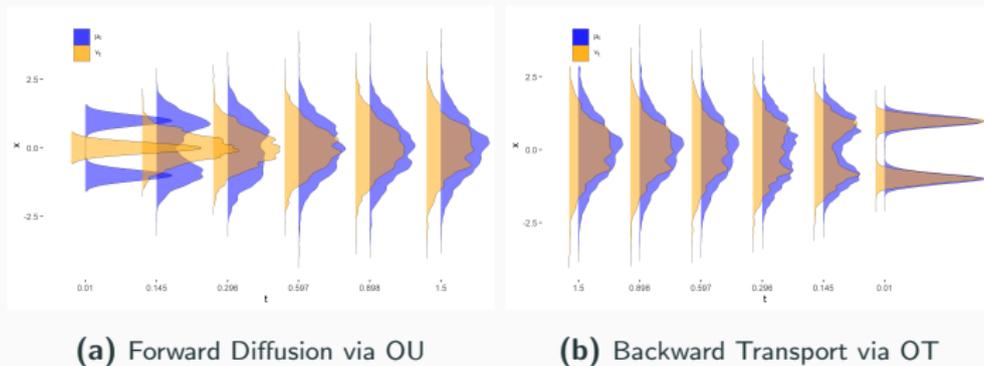


(c)

(a) $s_r(u)$ Low SNR. (b) $s_r(u)$ High SNR. (c) $\zeta_M^*(t)$.

non-log-concave example 1: two-point mass

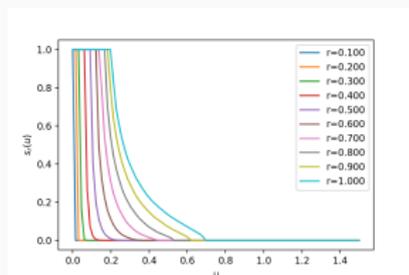
$$X_0 \sim \frac{1}{2}\delta_{-\mu} + \frac{1}{2}\delta_{\mu} \text{ for some } \mu > 0$$



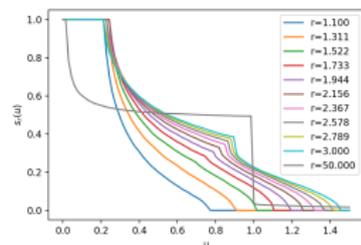
(a) Forward diffusion initialized at $\mu_0 = 0.5\delta_1 + 0.5\delta_{-1}$, $\nu_0 = \delta_0$. $T = 100$, $\eta = 0.01$. (b) Backward transport initialized at μ_T, ν_T , and applying the backward OT map $\mathbf{b}^{\mu, \eta}$.

non-log-concave example 2: mixture of point mass and Gaussian

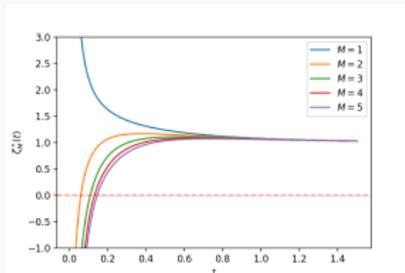
$$X_0 \sim \frac{1}{2}\delta_0 + \frac{1}{2}\mathcal{N}(0, 1)$$



(a)



(b)

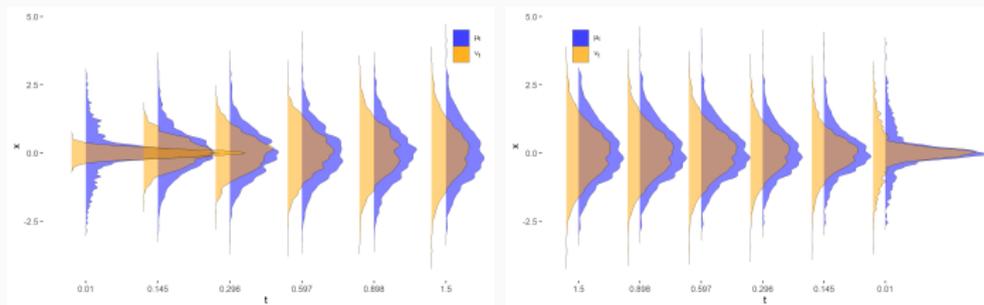


(c)

(a) $s_r(u)$ Low SNR. (b) $s_r(u)$ High SNR. (c) $\zeta_M^*(t)$.

non-log-concave example 2: mixture of point mass and Gaussian

$$X_0 \sim \frac{1}{2}\delta_0 + \frac{1}{2}\mathcal{N}(0, 1)$$



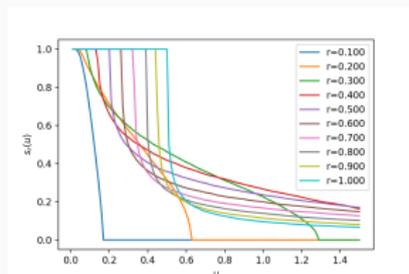
(a) Forward Diffusion via OU

(b) Backward Transport via OT

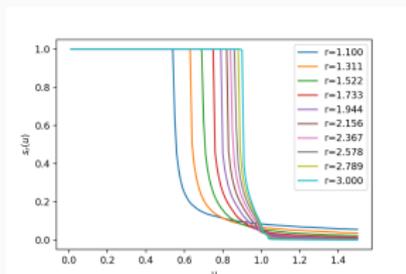
(a) Forward diffusion initialized at $\mu_0 = 0.5\delta_0 + 0.5\mathcal{N}(0, 1)$, $\nu_0 = \delta_0$. $T = 100$, $\eta = 0.01$. (b) Backward transport initialized at μ_T, ν_T , and applying the backward OT map $\mathbf{b}^{\mu, \eta}$.

non-log-concave example 3: heterogeneous Gaussian mixture

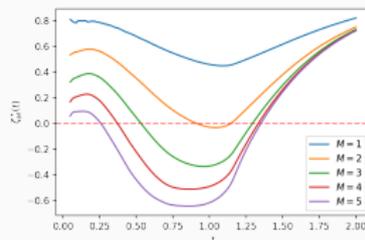
$$X_0 \sim \sum_{i=1}^m p_i \cdot \mathcal{N}(\mu_i, \sigma_i^2)$$



(a)



(b)

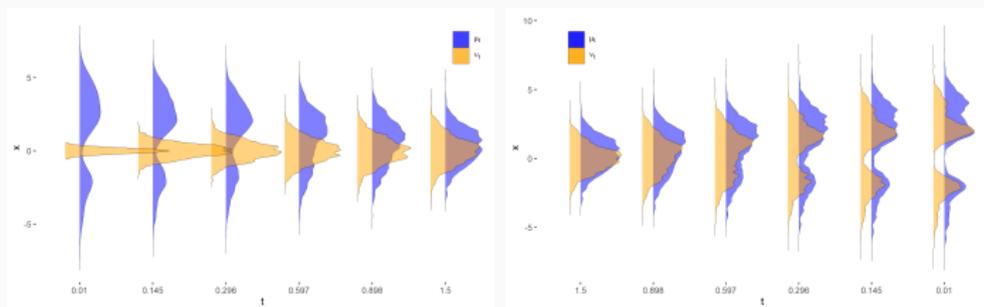


(c)

(a) $s_r(u)$ Low SNR. (b) $s_r(u)$ High SNR. (c) $\zeta_M^*(t)$.

non-log-concave example 3: heterogeneous Gaussian mixture

$$X_0 \sim \sum_{i=1}^m p_i \cdot \mathcal{N}(\mu_i, \sigma_i^2)$$



(a) Forward Diffusion via OU

(b) Backward Transport via OT

(a) Forward diffusion initialized at

$\mu_0 = 0.1\mathcal{N}(-4, 1) + 0.2\mathcal{N}(-2, 0.5) + 0.4\mathcal{N}(2, 0.5) + 0.3\mathcal{N}(4, 1)$, $\nu_0 = \delta_0$. $T = 100$, $\eta = 0.01$.

(b) Backward transport initialized at μ_T, ν_T , and applying the backward OT map $\mathbf{b}^{\mu, \eta}$.