An optimal control perspective on diffusion-based generative modeling

Julius Berner

Faculty of Mathematics. University of Vienna

December, 2022









Collaborators



Lorenz Richter

Zuse Institute Berlin, dida Datenschmiede GmbH



Karen Ullrich Meta Al

Generative modeling

Task

Sample from a complex (high-dimensional, multimodal) distribution \mathcal{D} .

Generative modeling

Task

Sample from a complex (high-dimensional, multimodal) distribution \mathcal{D} .

 $\ensuremath{\mathcal{D}}$ can be given in the form of:

1. samples $Y^{(i)} \sim D$ (images, text, sound, ...).



Generative modeling

Task

Sample from a complex (high-dimensional, multimodal) distribution \mathcal{D} .

 $\ensuremath{\mathcal{D}}$ can be given in the form of:

1. samples $Y^{(i)} \sim D$ (images, text, sound, ...).



2. an (unnormalized) **density** (e.g., in Bayesian statistics, computational physics and chemistry).



Fig. 2: https://en.wikipedia.org/wiki/File:Bimodal-bivariate-small.png Fig. 1: https://en.m.wikipedia.org/wiki/File:Cat_poster_1.jpg

Overview of generative models



Fig.: https://lilianweng.github.io/posts/2021-07-11-diffusion-models

Overview of generative models



This talk: **Diffusion models** (originated from (denoising) diffusion probabilistic modeling and score matching with Langevin dynamics).

J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In International Conference on Machine Learning, pages 2256-2265. PMLR, 2015

Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. Advances in Neural Information Processing Systems, 32, 2019

J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 202

Fig.: https://lilianweng.github.io/posts/2021-07-11-diffusion-models

Diffusion models

State-of-the-art models for a variety of applications where **samples** are given.



Sampling conditioned on the text prompt "a photograph of an astronaut riding a horse" using the stable diffusion model.

R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models, 2021

Julius Berner

Diffusion models

State-of-the-art models for a variety of applications where samples are given.
 Can be understood via stochastic differential equations (SDEs).



Sampling conditioned on the text prompt "a photograph of an astronaut riding a horse" using the stable diffusion model.



Fig.: https://yang-song.net/blog/2021/score

Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In International Conference on Learning Representations, 2020

R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models, 2021





R. Bellman. Dynamic programming. Princeton University Press, 1957

W. Fleming and R. Rishel. Deterministic and Stochastic Optimal Control. Applications of mathematics. Springer, 1975

A. Bensoussan. Perturbation methods in optimal control. Wiley, 1988





SDE-based generative modeling

Model: Stochastic differential equation (SDE)

d-dim. Brownian motion



drift

diffusion

SDE-based generative modeling

Model: Stochastic differential equation (SDE)

0.2

d-dim. Brownian motion



Simulation: For instance, via Euler-Maruyama scheme $\hat{X}_{t_k} \approx X_{t_k}$:

0.8



SDE-based generative modeling

Goal: Learn the drift $\bar{\mu}$ such that X_T approximates a given distribution \mathcal{D} .



Time-reversed SDE

Goal: $X_T \sim \mathcal{D}$ with generative SDE

 $\mathrm{d}X_s = \bar{\mu}(X_s,s)\,\mathrm{d}s + \bar{\sigma}(s)\,\mathrm{d}B_s.$

E. Nelson. Dynamical theories of Brownian motion. Press, Princeton, NJ, 1967

B. D. Anderson. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3):313-326, 1982

U. G. Haussmann and E. Pardoux. Time reversal of diffusions. The Annals of Probability, pages 1188-1205, 1986

Julius Berner

Time-reversed SDE

Goal: $X_T \sim \mathcal{D}$ with generative SDE $\mathrm{d}X_s = \tilde{\mu}(X_s, s) \,\mathrm{d}s + \tilde{\sigma}(s) \,\mathrm{d}B_s.$

"Solution": Define the inference SDE $dY_s = f(Y_s, s) ds + \sigma(s) dB_s, \ Y_0 \sim D,$ and set $X_0 \sim Y_T$ and $\mu \coloneqq \sigma \sigma^\top \nabla \log p_Y - f.$ $\Rightarrow p_X \stackrel{\text{a.e.}}{=} \bar{p}_Y$, i.e., $X_t \sim Y_{T-t}$ for $t \in [0, T].$



Notation: $\bar{\mu}(x, t) \coloneqq \mu(x, T - t)$ (time-reversal)

E. Nelson. Dynamical theories of Brownian motion. Press, Princeton, NJ, 1967

B. D. Anderson. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3):313-326, 1982

U. G. Haussmann and E. Pardoux. Time reversal of diffusions. The Annals of Probability, pages 1188-1205, 1986

Time-reversed SDE

Goal: $X_T \sim \mathcal{D}$ with generative SDE $\mathrm{d}X_s = \overline{\mu}(X_s, s) \,\mathrm{d}s + \overline{\sigma}(s) \,\mathrm{d}B_s.$

"Solution": Define the inference SDE $dY_s = f(Y_s, s) ds + \sigma(s) dB_s, \ Y_0 \sim D,$ and set $X_0 \sim Y_T$ and $\mu \coloneqq \sigma \sigma^\top \nabla \log p_Y - f.$ $\Rightarrow p_X \stackrel{\text{a.e.}}{=} \bar{p}_Y$, i.e., $X_t \sim Y_{T-t}$ for $t \in [0, T].$



Problem: How to learn the score $\nabla \log p_Y$ and how to sample from $X_0 \sim Y_T$?

E. Nelson. Dynamical theories of Brownian motion. Press, Princeton, NJ, 1967

B. D. Anderson. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3):313-326, 1982

U. G. Haussmann and E. Pardoux. Time reversal of diffusions. The Annals of Probability, pages 1188-1205, 1986

Julius Berner













Julius Berner

Hamilton-Jacobi-Bellman equation for time-reversed log-density

Define $V := -\log \bar{p}_X$ (Hopf-Cole transform), where $\bar{p}_X(x, t) := p_X(x, T - t)$ denotes the time-reversed density of the generative SDE X. Then

$$\partial_t V = -\operatorname{Tr}\left(D\nabla^2 V\right) + \mu \cdot \nabla V - \operatorname{div}(\mu) + \frac{1}{2} \left\|\sigma^\top \nabla V\right\|^2$$

Julius Berner

Deriving the ELBO



Deriving the ELBO



C.-W. Huang, J. H. Lim, and A. C. Courville. A variational perspective on diffusion-based generative models and score matching. Advances in Neural Information Processing Systems, 34, 2021

Verification theorem

Let V be the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D\nabla^2 V\right) - \mu \cdot \nabla V - h + \frac{1}{2} \|\sigma^\top \nabla V\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d} Y^{u}_{s} = (\sigma u - \mu) (Y^{u}_{s}, s) \, \mathrm{d} s + \sigma(s) \, \mathrm{d} B_{s},$$

be the controlled SDEs. Then it holds that

$$\mathbb{E}\left[V(Y_0^u,0)\right] = \min_{u \in \mathcal{U}} \mathbb{E}\left[\int_0^T \left(\frac{1}{2} \|u\|^2 + h\right) (Y_s^u,s) \, \mathrm{d}s + V(Y_T^u,T)\right].$$

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

Verification theorem

Let V be the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D\nabla^2 V\right) - \mu \cdot \nabla V - h + \frac{1}{2} \|\sigma^\top \nabla V\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d} Y^{u}_{s} = (\sigma u - \mu) (Y^{u}_{s}, s) \, \mathrm{d} s + \sigma(s) \, \mathrm{d} B_{s},$$

be the controlled SDEs. Then it holds that

$$\mathbb{E}\left[V(Y_0^u, 0)\right] = \min_{\substack{\psi \in \mathcal{U} \\ |\mathbf{t}\hat{o}'s \text{ lemma}}} \mathbb{E}\left[\underbrace{\int_0^T \left(\frac{1}{2} \|u\|^2 + h\right)(Y_s^u, s) \, \mathrm{d}s}_{\text{running costs}} + \underbrace{V(Y_T^u, T)}_{\text{terminal cost}}\right]$$

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

Verification theorem

Let V be the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D\nabla^2 V\right) - \mu \cdot \nabla V - h + \frac{1}{2} \|\sigma^\top \nabla V\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d}Y_{s}^{u} = (\sigma u - \mu)(Y_{s}^{u}, s)\mathrm{d}s + \sigma(s)\mathrm{d}B_{s},$$

be the controlled SDEs. Then it holds that

$$\mathbb{E}\left[V(Y_0^u,0)\right] = \min_{u \in \mathcal{U}} \mathbb{E}\left[\int_0^T \left(\frac{1}{2} \|u\|^2 + h\right) (Y_s^u,s) \, \mathrm{d}s + V(Y_T^u,T)\right]$$

where the unique minimum is attained by $u^* := -\sigma^\top \nabla V$.

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

Verification theorem

Let V be the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D\nabla^2 V\right) - \mu \cdot \nabla V - h + \frac{1}{2} \|\sigma^\top \nabla V\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d}Y_{s}^{u} = (\sigma u - \mu)(Y_{s}^{u}, s) \mathrm{d}s + \sigma(s) \mathrm{d}B_{s},$$

be the controlled SDEs. Then it holds that

$$\mathbb{E}\left[V(Y_0^u, 0)\right] = \mathbb{E}\left[\int_0^T \left(\frac{1}{2} \|u\|^2 + h\right) \left(Y_s^u, s\right) \mathrm{d}s + V(Y_T^u, T)\right] - \underbrace{D_{\mathsf{KL}}(\mathbb{P}_{Y^u}|\mathbb{P}_{Y^{u^*}})}_{\mathsf{KL}-\mathsf{divergence on path space}},$$

where $u^* \coloneqq -\sigma^\top \nabla V$.

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

Verification theorem

Let V be the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D\nabla^2 V\right) - \mu \cdot \nabla V - \operatorname{div}(\mu) + \frac{1}{2} \left\|\sigma^\top \nabla V\right\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d}Y_{s}^{u} = (\sigma u - \mu)(Y_{s}^{u}, s) \mathrm{d}s + \sigma(s) \mathrm{d}B_{s},$$

be the controlled SDEs. Then it holds that

$$\mathbb{E}\left[V(Y_0^u, 0)\right] = \mathbb{E}\left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\mu)\right)(Y_s^u, s) \,\mathrm{d}s + V(Y_T^u, T)\right] - D_{\mathsf{KL}}(\mathbb{P}_{Y^u}|\mathbb{P}_{Y^{u^*}}),$$

where $u^* \coloneqq -\sigma^\top \nabla V.$

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

ELBO

Let $V := -\log \overline{p}_X$ (time-reversed log-density), i.e., the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D
abla^2 V
ight) - \mu\cdot
abla V - \operatorname{div}(\mu) + rac{1}{2} \left\|\sigma^{ op}
abla V
ight\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d} Y^u_s = (\sigma u - \mu) (Y^u_s, s) \mathrm{d} s + \sigma(s) \mathrm{d} B_s,$$

be the controlled SDEs. Then it holds that

$$\mathbb{E}\left[\log p_{X_{T}}(Y_{0}^{u})\right] = \mathbb{E}\left[-\int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\mu)\right)(Y_{s}^{u}, s) \,\mathrm{d}s + \log p_{X_{0}}(Y_{T}^{u})\right] + D_{\mathsf{KL}}(\mathbb{P}_{Y^{u}}|\mathbb{P}_{Y^{u^{*}}}),$$

where $u^{*} \coloneqq \sigma^{\top} \nabla \log p_{X}$ (scaled score).

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

ELBO

Let $V := -\log \overline{p}_X$ (time-reversed log-density), i.e., the solution to the HJB equation

$$\partial_t V = -\operatorname{Tr}\left(D
abla^2 V
ight) - \mu\cdot
abla V - \operatorname{div}(\mu) + rac{1}{2} \left\|\sigma^{ op}
abla V
ight\|^2,$$

let $\ensuremath{\mathcal{U}}$ be a suitable space of controls, and let

$$\mathrm{d}Y^u_s = (\sigma u - \mu)(Y^u_s, s)\,\mathrm{d}s + \sigma(s)\,\mathrm{d}B_s,$$

be the controlled SDEs. Then it holds that

$$\underbrace{\mathbb{E}\left[\log p_{X_{T}}(Y_{0}^{u})\right]}_{\text{log-likelihood}} = \mathbb{E}\left[-\int_{0}^{T}\left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\mu)\right)(Y_{s}^{u}, s) \,\mathrm{d}s + \log p_{X_{0}}(Y_{T}^{u})\right] + \underbrace{D_{\mathrm{KL}}(\mathbb{P}_{Y^{u}}|\mathbb{P}_{Y^{u^{*}}})}_{\text{variational gap}}$$

where $u^* \coloneqq \sigma^\top \nabla \log \overline{p}_X$ (scaled score).

M. Pavon. Stochastic control and nonequilibrium thermodynamical systems. Applied Mathematics and Optimization, 19(1):187-202, 1989

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

Deriving the ELBO

ELBO:
$$\mathbb{E}\left[\log p_{X_T}(Y_0^u)\right] \ge \mathbb{E}\left[\log p_{X_0}(Y_T^u) - \int_0^T \left(\frac{1}{2}\|u\|^2 + \operatorname{div}(\mu)\right)(Y_s^u, s) \,\mathrm{d}s\right]$$

SDEs: $\mathrm{d}Y_s^u = (\sigma u - \mu)(Y_s^u, s) \,\mathrm{d}s + \sigma(s) \,\mathrm{d}B_s, \quad \mathrm{d}X_s = \overline{\mu}(X_s, s) \,\mathrm{d}s + \overline{\sigma}(s) \,\mathrm{d}B_s$
Deriving the ELBO

ELBO:
$$\mathbb{E}\left[\log p_{X_T}(Y_0^u)\right] \ge \mathbb{E}\left[\log p_{X_0}(Y_T^u) - \int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\mu)\right)(Y_s^u, s) \,\mathrm{d}s\right]$$

SDEs: $\mathrm{d}Y_s^u = (\sigma u - \mu)(Y_s^u, s) \,\mathrm{d}s + \sigma(s) \,\mathrm{d}B_s, \quad \mathrm{d}X_s = \tilde{\mu}(X_s, s) \,\mathrm{d}s + \tilde{\sigma}(s) \,\mathrm{d}B_s$

Problem: We want to control the generative SDE X.

Deriving the ELBO

ELBO:
$$\mathbb{E}\left[\log p_{X_T}(Y_0^u)\right] \ge \mathbb{E}\left[\log p_{X_0}(Y_T^u) - \int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\mu)\right)(Y_s^u, s) \,\mathrm{d}s\right]$$

SDEs: $dY_s^u = (\sigma u - \mu) (Y_s^u, s) ds + \sigma(s) dB_s, \quad dX_s = \overline{\mu}(X_s, s) ds + \overline{\sigma}(s) dB_s$

- **Problem:** We want to control the generative SDE X.
- \Rightarrow **Reparametrization:** $\mu \coloneqq \sigma u f$

$$\begin{aligned} & \textbf{ELBO:} \quad \mathbb{E}\left[\log p_{X_{T}^{u}}(Y_{0})\right] \geq \mathbb{E}\left[\log p_{X_{0}^{u}}(Y_{T}) - \int_{0}^{T} \left(\frac{1}{2} \|u\|^{2} + \operatorname{div}(\sigma u - f)\right)(Y_{s}, s) \, \mathrm{d}s\right] \\ & \textbf{SDEs:} \quad \mathrm{d}Y_{s} = f(Y_{s}, s) \, \mathrm{d}s + \sigma(s) \, \mathrm{d}B_{s}, \qquad \mathrm{d}X_{s}^{u} = \left(\overleftarrow{\sigma}\overleftarrow{u} - \overleftarrow{f}\right)(X_{s}^{u}, s) \, \mathrm{d}s + \overleftarrow{\sigma}(s) \, \mathrm{d}B_{s} \end{aligned}$$

C.-W. Huang, J. H. Lim, and A. C. Courville. A variational perspective on diffusion-based generative models and score matching. Advances in Neural Information Processing Systems, 34, 2021

Score-based generative modeling

$$\begin{array}{l} \textbf{ELBO:} \quad \overbrace{\mathbb{E}\left[\log p_{X_{T}^{u}}(Y_{0})\right]}^{\text{log-likelihood}} \geq \mathbb{E}\left[\overbrace{\log p_{X_{0}^{u}}(Y_{T})}^{\text{prior loss}} - \int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\sigma u - f)\right)(Y_{s}^{u}, s) \,\mathrm{d}s\right] \\ \textbf{SDEs:} \quad \underbrace{\mathrm{d}Y_{s} = f(Y_{s}, s) \,\mathrm{d}s + \sigma(s) \,\mathrm{d}B_{s}, \ Y_{0} \sim \mathcal{D}}_{\text{inference SDE}}, \ \underbrace{\mathrm{d}X_{s}^{u} = \left(\overline{\sigma}\,\overline{u} - \overline{f}\right)(X_{s}^{u}, s) \,\mathrm{d}s + \overline{\sigma}(s) \,\mathrm{d}B_{s}}_{\text{generative SDE}} \end{aligned}$$

Score-based generative modeling

$$\begin{array}{l} \textbf{ELBO:} \quad \overbrace{\mathbb{E}\left[\log p_{X_{T}^{u}}(Y_{0})\right]}^{\text{log-likelihood}} \geq \mathbb{E}\left[\overbrace{\log p_{X_{0}^{u}}(Y_{T})}^{\text{prior loss}} - \int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\sigma u - f)\right)(Y_{s}^{u}, s) \,\mathrm{d}s\right] \\ \textbf{SDEs:} \quad \underbrace{\mathrm{d}Y_{s} = f(Y_{s}, s) \,\mathrm{d}s + \sigma(s) \,\mathrm{d}B_{s}, \ Y_{0} \sim \mathcal{D}}_{\text{inference SDE}}, \ \underbrace{\mathrm{d}X_{s}^{u} = \left(\overline{\sigma}\,\overline{u} - \overline{f}\right)(X_{s}^{u}, s) \,\mathrm{d}s + \overline{\sigma}(s) \,\mathrm{d}B_{s}}_{\text{generative SDE}} \end{aligned}$$

Time-reversal: If $X_0^u \sim Y_T$, then $u = u^* = \sigma^\top \nabla \log p_Y$ and $X_T^u \sim \mathcal{D}$.



Score-based generative modeling

$$\begin{array}{l} \textbf{ELBO:} \quad \overbrace{\mathbb{E}\left[\log p_{X_{T}^{u}}(Y_{0})\right]}^{\text{log-likelihood}} \geq \mathbb{E}\left[\overbrace{\log p_{X_{0}^{u}}(Y_{T})}^{\text{prior loss}} - \int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\sigma u - f)\right)(Y_{s}^{u}, s) \,\mathrm{d}s\right] \\ \textbf{SDEs:} \quad \underbrace{\mathrm{d}Y_{s} = f(Y_{s}, s) \,\mathrm{d}s + \sigma(s) \,\mathrm{d}B_{s}, \ Y_{0} \sim \mathcal{D}}_{\text{inference SDE}}, \ \underbrace{\mathrm{d}X_{s}^{u} = \left(\overline{\sigma}\,\overline{u} - \overline{f}\right)(X_{s}^{u}, s) \,\mathrm{d}s + \overline{\sigma}(s) \,\mathrm{d}B_{s}}_{\text{generative SDE}} \end{aligned}$$

Time-reversal: If $X_0^u \sim Y_T$, then $u = u^* = \sigma^\top \nabla \log p_Y$ and $X_T^u \sim \mathcal{D}$.



Questions: 1. Can the ELBO be simplified? 2. How to sample from $X_0^u \sim Y_T$?

Julius Berner

Denoising score matching

1. Rewrite the divergence and use a MC approximation, $\tau \sim \mathcal{U}([0, T])$, of the time-integral:

Denoising score matching

1. Rewrite the divergence and use a MC approximation, $\tau \sim \mathcal{U}([0, T])$, of the time-integral:

ELBO:
$$\mathbb{E}\left[\log p_{X_T^u}(Y_0)\right] \geq \mathbb{E}\left[\log p_{X_0^u}(Y_T) - \int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\sigma u - f)\right)(Y_s, s) \, \mathrm{d}s\right]$$

P. Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661-1674, 2011

C.-W. Huang, J. H. Lim, and A. C. Courville. A variational perspective on diffusion-based generative models and score matching. Advances in Neural Information Processing Systems, 34, 2021

Julius Berner

An optimal control perspective on diffusion-based generative modeling

16 / 30

Denoising score matching

1. Rewrite the divergence and use a MC approximation, $\tau \sim \mathcal{U}([0, T])$, of the time-integral:

Divergence theorem:
$$u \cdot \sigma^{\top} \nabla \log p_{Y_{s}|Y_{0}}$$
 (in expectation)
ELBO: $\mathbb{E} \left[\log p_{X_{T}^{u}}(Y_{0}) \right] \geq \mathbb{E} \left[\log p_{X_{0}^{u}}(Y_{T}) - \int_{0}^{T} \left(\frac{1}{2} \|u\|^{2} + \operatorname{div}(\sigma u) - \operatorname{div}(f) \right)(Y_{s}, s) \, \mathrm{d}s \right]$

$$= \frac{T}{2} \underbrace{\mathbb{E} \left[\left\| u(Y_{\tau}, \tau) - \sigma^{\top}(\tau) \nabla \log p_{Y_{\tau}|Y_{0}}(Y_{\tau}|Y_{0}) \right\|^{2} \right]}_{\operatorname{denoising score matching}} + \operatorname{const.}$$

P. Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661-1674, 2011

C.-W. Huang, J. H. Lim, and A. C. Courville. A variational perspective on diffusion-based generative models and score matching. Advances in Neural Information Processing Systems, 34, 2021

Ornstein–Uhlenbeck process

2. Construct drift f and diffusion σ of Y such that $Y_T \approx \mathcal{N}(0, I)$:

Ornstein-Uhlenbeck process

2. Construct drift f and diffusion σ of Y such that $Y_T \approx \mathcal{N}(0, I)$:

SDEs:
$$dX_s^u = (\bar{\sigma}\bar{u} - \bar{f}) (X_s^u, s) ds + \bar{\sigma}(s) dB_s, \quad X_0 \sim \mathcal{N}(0, I)$$

 $dY_s = f(Y_s, s) ds + \sigma(s) dB_s, \quad Y_0 \sim \mathcal{D}$
 $\uparrow \qquad \uparrow$
 $= b(s)Y_s \quad \text{diagonal matrix}$



Ornstein-Uhlenbeck process

2. Construct drift f and diffusion σ of Y such that $Y_T \approx \mathcal{N}(0, I)$:

SDEs:
$$dX_s^u = (\bar{\sigma}\bar{u} - \bar{f}) (X_s^u, s) ds + \bar{\sigma}(s) dB_s, \quad X_0 \sim \mathcal{N}(0, I)$$

 $dY_s = f(Y_s, s) ds + \sigma(s) dB_s, \quad Y_0 \sim \mathcal{D}$
 $\uparrow \qquad \uparrow$
 $= b(s)Y_s \qquad \text{diagonal matrix}$



Example: Variance-preserving SDE (Ornstein-Uhlenbeck process) given by $b = -\frac{1}{2}\sigma^2$ such that

$$Y_t|Y_0 \sim \mathcal{N}\left(e^{-lpha(t)}Y_0, 1-e^{-2lpha(t)}
ight),$$

where $\alpha(t) \coloneqq \frac{1}{2} \int_0^t \sigma^2(s) ds$.

Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In International Conference on Learning Representations, 2020

Summary

variance-preserving SDE with $Y_T \approx \mathcal{N}(0, I)$ for sufficiently large $\alpha(T) := \frac{1}{2} \int_0^T \sigma^2(s) ds$

SDEs:
$$\underbrace{\mathrm{d}Y_s = -\frac{1}{2}\sigma^2(s)Y_s\,\mathrm{d}s + \sigma(s)\,\mathrm{d}B_s}_{\mathrm{d}S}, \quad Y_0 \sim \mathcal{D}$$
$$\underbrace{\mathrm{d}X_s^u = \left(\bar{\sigma}\bar{u} - \bar{f}\right)\left(X_s^u, s\right)\mathrm{d}s + \bar{\sigma}(s)\,\mathrm{d}B_s}_{\mathrm{generative SDE}}, \quad X_0 \sim \mathcal{N}(0, \mathrm{I})$$



Summary



Summary

variance-preserving SDE with $Y_T \approx \mathcal{N}(0, I)$ for sufficiently large $\alpha(T) := \frac{1}{2} \int_0^T \sigma^2(s) ds$ $\mathrm{d}Y_s = -\frac{1}{2}\sigma^2(s)Y_s\,\mathrm{d}s + \sigma(s)\,\mathrm{d}B_s, \quad Y_0 \sim \mathcal{D}$ SDEs: $\mathrm{d} X^u_s = \left(\bar{\sigma} \bar{u} - \bar{f}\right) (X^u_s, s) \, \mathrm{d} s + \bar{\sigma}(s) \, \mathrm{d} B_s, \quad X_0 \sim \mathcal{N}(0, \mathrm{I})$ generative SDE log-likelihood denoising score matching $\textbf{ELBO:} \quad \mathbb{E}\left[\log p_{X_T^u}(Y_0)\right] \geq \mathbb{E}\left[\left\|u(Y_{\tau},\tau) - \sigma(\tau)\nabla \log p_{Y_{\tau}|Y_0}(Y_{\tau}|Y_0)\right\|^2\right], \quad \tau \sim \mathcal{U}([0,T])$ $= \mathcal{N}\left(e^{-\alpha(t)}Y_0, 1-e^{-2\alpha(t)}\right)$

 \Rightarrow We arrived at (a variant of) the objective used to train diffusion models in practice!

Julius Berner



Setting: Data distribution \mathcal{D} admits the density $\frac{\rho}{\mathcal{Z}}$, where ρ is known but the normalizing constant $\mathcal{Z} := \int_{\mathbb{R}^d} \rho(x) \, dx$ might be intractable.



Setting: Data distribution \mathcal{D} admits the density $\frac{\rho}{\mathcal{Z}}$, where ρ is known but the normalizing constant $\mathcal{Z} := \int_{\mathbb{R}^d} \rho(x) \, dx$ might be intractable.



Problem: Cannot use the previous ELBO as we cannot simulate Y starting at $Y_0 \sim D$:

$$\mathbb{E}\left[\log p_{X_{\tau}^{u}}(Y_{0})\right] \geq \mathbb{E}\left[\left\|u(Y_{\tau},\tau)-\sigma(\tau)\nabla\log p_{Y_{\tau}|Y_{0}}(Y_{\tau}|Y_{0})\right\|^{2}\right].$$

Idea: Use verification theorem on HJB equation for $V := -\log \bar{p}_Y$ (instead of $-\log \bar{p}_X$):

$$\partial_t V = -\operatorname{Tr}\left(\overline{D}\nabla^2 V\right) + \overline{f} \cdot \nabla V - \operatorname{div}(\overline{f}) + \frac{1}{2} \|\overline{\sigma}^\top \nabla V\|^2.$$

Idea: Use verification theorem on HJB equation for $V := -\log \bar{p}_Y$ (instead of $-\log \bar{p}_X$):

$$\partial_t V = -\operatorname{Tr}\left(\overline{D}\nabla^2 V\right) + \overline{f} \cdot \nabla V - \operatorname{div}(\overline{f}) + \frac{1}{2} \|\overline{\sigma}^\top \nabla V\|^2.$$

Verification theorem

Define the controlled SDEs

$$\mathrm{d} X^u_s = \left(\overleftarrow{\sigma} u - \overleftarrow{f} \right) (X^u_s, s) \, \mathrm{d} s + \overleftarrow{\sigma}(s) \, \mathrm{d} B_s.$$

Then it holds that

$$\mathbb{E}\left[\log p_{Y_{\mathcal{T}}}(X_0^u)\right] = \mathbb{E}\left[\log p_{Y_0}(X_{\mathcal{T}}^u) - \int_0^{\mathcal{T}} \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f})\right) (X_s^u, s) \,\mathrm{d}s\right] + D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{X^{u^*}}).$$

where $u^* = \overleftarrow{\sigma}^\top \nabla \log \overleftarrow{p}_Y$.

Julius Berner

Verification theorem:

$$\mathbb{E}\left[\log p_{Y_{\mathcal{T}}}(X_0^u)\right] = \mathbb{E}\left[\log p_{Y_0}(X_T^u) - \int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f})\right)(X_s^u, s) \,\mathrm{d}s\right] + D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{X^{u^*}})$$

Verification theorem:

$$\mathbb{E}\left[\log p_{Y_{T}}(X_{0}^{u})\right] = \mathbb{E}\left[\log p_{Y_{0}}(X_{T}^{u}) - \int_{0}^{T} \left(\frac{1}{2} \|u\|^{2} + \operatorname{div}(\tilde{f})\right)(X_{s}^{u}, s) \,\mathrm{d}s\right] + D_{\mathsf{KL}}(\mathbb{P}_{X^{u}}|\mathbb{P}_{X^{u^{*}}})$$

Using $p_{Y_0} = \frac{\rho}{Z}$, yields the following reverse KL objective:

$$D_{\mathsf{KL}}(\mathbb{P}_{X^{u}}|\mathbb{P}_{X^{u^{*}}}) = \mathbb{E}\left[\int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \mathsf{div}(\tilde{f})\right)(X^{u}_{s}, s) \,\mathrm{d}s + \log \frac{p_{Y_{T}}(X^{u}_{0})}{\rho(X^{u}_{T})}\right] + \log \mathcal{Z}$$

Verification theorem:

$$\mathbb{E}\left[\log p_{Y_{T}}(X_{0}^{u})\right] = \mathbb{E}\left[\log p_{Y_{0}}(X_{T}^{u}) - \int_{0}^{T} \left(\frac{1}{2} \|u\|^{2} + \operatorname{div}(\tilde{f})\right)(X_{s}^{u}, s) \,\mathrm{d}s\right] + D_{\mathsf{KL}}(\mathbb{P}_{X^{u}}|\mathbb{P}_{X^{u^{*}}})$$

Using $p_{Y_0} = \frac{\rho}{Z}$, yields the following **reverse KL objective**:

$$D_{\mathsf{KL}}(\mathbb{P}_{X^{u}}|\mathbb{P}_{X^{u^{*}}}) = \mathbb{E}\left[\int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\bar{f})\right)(X_{s}^{u}, s) \,\mathrm{d}s + \log \frac{\overbrace{p_{Y_{T}}(X_{0}^{u})}{\rho(X_{T}^{u})}}{\rho(X_{T}^{u})}\right] + \underbrace{\log \mathcal{Z}}_{\operatorname{const.}}$$

Verification theorem:

$$\mathbb{E}\left[\log p_{Y_{T}}(X_{0}^{u})\right] = \mathbb{E}\left[\log p_{Y_{0}}(X_{T}^{u}) - \int_{0}^{T} \left(\frac{1}{2} \|u\|^{2} + \operatorname{div}(\tilde{f})\right)(X_{s}^{u}, s) \,\mathrm{d}s\right] + D_{\mathsf{KL}}(\mathbb{P}_{X^{u}}|\mathbb{P}_{X^{u^{*}}})$$

Using $p_{Y_0} = \frac{\rho}{Z}$, yields the following reverse KL objective:

$$D_{\mathsf{KL}}(\mathbb{P}_{X^{u}}|\mathbb{P}_{X^{u^{*}}}) = \mathbb{E}\left[\int_{0}^{T} \left(\frac{1}{2}\|u\|^{2} + \operatorname{div}(\overline{f})\right)(X_{s}^{u}, s) \,\mathrm{d}s + \log\frac{\overbrace{\rho_{Y_{T}}(X_{0}^{u})}{\rho(X_{T}^{u})}}{\rho(X_{T}^{u})}\right] + \underbrace{\log \mathcal{Z}}_{\operatorname{const.}}$$

prior loss

Remark: $D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{X^{u^*}}) = D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{\bar{Y}}) - \overbrace{D_{\mathsf{KL}}(\mathbb{P}_{X^u_0}|\mathbb{P}_{Y_T})}^{\mathcal{D}}$ with reverse time SDE \bar{Y} .



⇒ Diffusion-based method for (unnormalized) densities: Time-Reversed **Di**ffusion **S**ampler

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\bar{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \\ \mathrm{d}X_s^u = \left(\bar{\sigma}u - \bar{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0^u \sim \mathcal{N}(0, \mathrm{I})$$

 \Rightarrow Diffusion-based method for (unnormalized) densities: Time-Reversed **Di**ffusion **S**ampler

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \\ \mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0^u \sim \mathcal{N}(0, \mathrm{I})$$

Cannot simplify the objective (as in denoising score matching) since the expectation is over the controlled process X^u.

 \Rightarrow Diffusion-based method for (unnormalized) densities: Time-Reversed **Di**ffusion **S**ampler

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \\ \mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0^u \sim \mathcal{N}(0, \mathrm{I})$$

- Cannot simplify the objective (as in denoising score matching) since the expectation is over the controlled process X^u.
- \Rightarrow Discretize objective, parametrize *u* as neural network, simulate whole trajectory X^{u} , and backpropagate (or use stochastic adjoint method).

P. Kidger, J. Foster, X. C. Li, and T. Lyons. Efficient and accurate gradients for neural SDEs. Advances in Neural Information Processing Systems, 34:18747-18761, 2021

 \Rightarrow Diffusion-based method for (unnormalized) densities: Time-Reversed **Di**ffusion **S**ampler

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \\ \mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0^u \sim \mathcal{N}(0, \mathrm{I})$$

- Cannot simplify the objective (as in denoising score matching) since the expectation is over the controlled process X^u.
- \Rightarrow Discretize objective, parametrize *u* as neural network, simulate whole trajectory X^{u} , and backpropagate (or use stochastic adjoint method).
- Initialize neural network as linear interpolation between (scaled) initial and terminal scores $\bar{\sigma}^{\top}\nabla \log \mathcal{N}(0, \mathbf{I}) \approx \bar{\sigma}^{\top}\nabla \log p_{Y_{T}}$ and $\bar{\sigma}^{\top}\nabla \log \rho = \bar{\sigma}^{\top}\nabla \log p_{Y_{0}}$.

P. Kidger, J. Foster, X. C. Li, and T. Lyons. Efficient and accurate gradients for neural SDEs. Advances in Neural Information Processing Systems, 34:18747–18761, 2021

Double-Well: Sampling from a double-well (DW) with 2⁵ modes in 20 dimensions:

$$\rho(x) = \exp\left(-\sum_{i=1}^{w}(x_i^2-\delta)^2 - \frac{1}{2}\sum_{i=w+1}^{d}x_i^2\right), \quad d = 20, \ \delta = 3, \ w = 5.$$



An optimal control perspective on diffusion-based generative modeling

24 / 30

We compare against one of the standard approaches (Path Integral Sampler) originating from the field of stochastic optimal control.

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313–329, 1991
B. Tzen and M. Raginsky. Theoretical guarantees for sampling and inference in generative models with latent diffusions. In *Conference on Learning Theory*, pages 3084–3114. PMLR, 2019
 Q. Zhang and Y. Chen. Path integral sampler: a stochastic control approach for sampling. In *Internet Conference on Learning Theory*, pages 3084–3114. PMLR, 2019
 L. Richter. Solving high-dimensional *PDEs*, approximation of path space measures and importance sampling of diffusions. PhD thesis, PHD thesis, PHD thesis, PHD thesis, PHD thesis, PHD thesis, SHD Cottbus-Senftenberg, 2021
 F. Vargas, A. Ovaianas, D. Fernades, M. Girolami, N. D. Lawrence, and N. Nisken. Bayesian learning via neural Schrödinger-Föllmer (How. Statistics and Computing, 33(1)-12-22, 2033

- We compare against one of the standard approaches (Path Integral Sampler) originating from the field of stochastic optimal control.
- Based on Schrödinger (Half-)Bridge problem $\min_{u \in \mathcal{U}} D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{X^0})$ s.t. $X_T^u \sim \mathcal{D}$, where

$$\mathrm{d} X^u_s = igl(\sigma u + \muigl) (X^u_s, s) \,\mathrm{d} s + \sigma(s) \,\mathrm{d} B_s, \quad X^u_0 \sim \delta_{x_0}.$$

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313–329, 1991
 B. Tzen and M. Raginsky. Theoretical guarantees for sampling and inference in generative models with latent diffusions. In *Conference on Learning Theory*, pages 3084–3114, PMLR, 2019
 Q. Zhang and Y. Chen. Path Integral sampler: a stochastic control approach for sampling. In *International Conference on Learning Representations*, 2022
 L. Richter. Solving high-dimensional *PDEs*, approximation of path space measures and importance sampling of diffusions. PhD thesis, BTU Cottbus Senftenberg, 2021
 F. Vargas, A. Ovaianas, D. Fernandes, M. Giralami, N. D. Lawrence, and N. Nisken. Bayesian learning via neural Schödinger–Föllmer flows. *Statistics and Computing*, 33(1):1–22, 2003

- We compare against one of the standard approaches (Path Integral Sampler) originating from the field of stochastic optimal control.
- Based on Schrödinger (Half-)Bridge problem $\min_{u \in U} D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{X^0})$ s.t. $X_T^u \sim \mathcal{D}$, where

$$\mathrm{d}X^u_s = (\sigma u + \mu)(X^u_s, s)\,\mathrm{d}s + \sigma(s)\,\mathrm{d}B_s, \quad X^u_0 \sim \delta_{x_0}.$$

Leads to the following objective:

PIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|u(X_s^u, s)\|^2 \, \mathrm{d}s + \log \frac{p_{X_T^0}(X_T^u)}{\rho(X_T^u)} \right]$$

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

B. Tzen and M. Raginsky. Theoretical guarantees for sampling and inference in generative models with latent diffusions. In Conference on Learning Theory, pages 3084–3114. PMLR, 2019

Q. Zhang and Y. Chen. Path integral sampler: a stochastic control approach for sampling. In International Conference on Learning Representations, 2022

L. Richter. Solving high-dimensional PDEs, approximation of path space measures and importance sampling of diffusions. PhD thesis, BTU Cottbus-Senftenberg, 2021

F. Vargas, A. Ovsianas, D. Fernandes, M. Girolami, N. D. Lawrence, and N. Nüsken. Bayesian learning via neural Schrödinger-Föllmer flows. Statistics and Computing, 33(1):1-22, 2023

Julius Berner

An optimal control perspective on diffusion-based generative modeling

25 / 30

- We compare against one of the standard approaches (Path Integral Sampler) originating from the field of stochastic optimal control.
- Based on Schrödinger (Half-)Bridge problem $\min_{u \in U} D_{\mathsf{KL}}(\mathbb{P}_{X^u}|\mathbb{P}_{X^0})$ s.t. $X_T^u \sim \mathcal{D}$, where

$$\mathrm{d} X^{u}_{s} = \big(\sigma u + \widehat{\mu}\,\big)(X^{u}_{s},s)\,\mathrm{d} s + \overbrace{\sigma(s)}^{=\frac{1}{\sqrt{\tau}}} \mathrm{d} B_{s}, \quad X^{u}_{0} \sim \delta_{x_{0}}.$$

Leads to the following objective:

PIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|u(X_s^u, s)\|^2 \, \mathrm{d}s + \log \frac{\widetilde{p_{X_T^0}(X_T^u)}}{\rho(X_T^u)} \right]$$

P. Dai Pra. A stochastic control approach to reciprocal diffusion processes. Applied mathematics and Optimization, 23(1):313-329, 1991

B. Tzen and M. Raginsky. Theoretical guarantees for sampling and inference in generative models with latent diffusions. In Conference on Learning Theory, pages 3084–3114. PMLR, 2019

Q. Zhang and Y. Chen. Path integral sampler: a stochastic control approach for sampling. In International Conference on Learning Representations, 2022

L. Richter. Solving high-dimensional PDEs, approximation of path space measures and importance sampling of diffusions. PhD thesis, BTU Cottbus-Senftenberg, 2021

F. Vargas, A. Ovsianas, D. Fernandes, M. Girolami, N. D. Lawrence, and N. Nüsken. Bayesian learning via neural Schrödinger-Föllmer flows. Statistics and Computing, 33(1):1-22, 2023

Julius Berner

An optimal control perspective on diffusion-based generative modeling

25 / 30

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right]^{-1} \approx \text{density of } Y_T$$
$$\mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0 \sim \mathcal{N}(0, \mathrm{I})$$

PIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|u(X_s^u, s)\|^2 \, \mathrm{d}s + \log \frac{\mathcal{N}(X_T^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \quad \text{density of } X_T^0$$
$$\mathrm{d}X_s^u = \frac{1}{\sqrt{T}} u(X_s^u, s) \, \mathrm{d}s + \frac{1}{\sqrt{T}} \, \mathrm{d}B_s, \quad X_0 \sim \delta_{x_0}$$

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right]^{-1} \approx \text{density of } Y_T$$
$$\mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0 \sim \mathcal{N}(0, \mathrm{I})$$

PIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|u(X_s^u, s)\|^2 \, \mathrm{d}s + \log \frac{\mathcal{N}(X_T^u; 0, \mathrm{I})}{\rho(X_T^u)} \right]^{-1} \text{ density of } X_T^0$$
$$\mathrm{d}X_s^u = \frac{1}{\sqrt{T}} u(X_s^u, s) \, \mathrm{d}s + \frac{1}{\sqrt{T}} \, \mathrm{d}B_s, \quad X_0 \sim \delta_{x_0}$$

 \Rightarrow DIS can use different SDEs and initial distributions $X \sim \mathcal{N}(0,
u^2 \mathrm{I})$

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right]^{-1} \approx \text{density of } Y_T$$
$$\mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0 \sim \mathcal{N}(0, \mathrm{I})$$

PIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|u(X_s^u, s)\|^2 \, \mathrm{d}s + \log \frac{\mathcal{N}(X_T^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \quad \text{density of } X_T^0$$
$$\mathrm{d}X_s^u = \frac{1}{\sqrt{T}} u(X_s^u, s) \, \mathrm{d}s + \frac{1}{\sqrt{T}} \, \mathrm{d}B_s, \quad X_0 \sim \delta_{x_0}$$

- \Rightarrow DIS can use different SDEs and initial distributions $X \sim \mathcal{N}(0,
 u^2 \mathrm{I})$
- \Rightarrow Allows to transfer techniques for diffusion models (noise schedules, probability flow ODE, exponential integrators, ...)

Q. Zhang and Y. Chen. Fast sampling of diffusion models with exponential integrator. arXiv preprint arXiv:2204.13902, 2022

DIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \left(\frac{1}{2} \|u\|^2 + \operatorname{div}(\tilde{f}) \right) (X_s^u, s) \, \mathrm{d}s + \log \frac{\mathcal{N}(X_0^u; 0, \mathrm{I})}{\rho(X_T^u)} \right]^{-\infty} \approx \text{density of } Y_T$$
$$\mathrm{d}X_s^u = \left(\bar{\sigma}u - \tilde{f} \right) (X_s^u, s) \, \mathrm{d}s + \bar{\sigma}(s) \, \mathrm{d}B_s, \quad X_0 \sim \mathcal{N}(0, \mathrm{I})$$

PIS:
$$\min_{u \in \mathcal{U}} \mathbb{E} \left[\int_0^T \frac{1}{2} \|u(X_s^u, s)\|^2 \, \mathrm{d}s + \log \frac{\mathcal{N}(X_T^u; 0, \mathrm{I})}{\rho(X_T^u)} \right] \quad \text{density of } X_T^0$$
$$\mathrm{d}X_s^u = \frac{1}{\sqrt{T}} u(X_s^u, s) \, \mathrm{d}s + \frac{1}{\sqrt{T}} \, \mathrm{d}B_s, \quad X_0 \sim \delta_{x_0}$$

- $\Rightarrow\,$ DIS can use different SDEs and initial distributions $X\sim\mathcal{N}(0,
 u^2\mathrm{I})$
- ⇒ Allows to transfer techniques for diffusion models (noise schedules, probability flow ODE, exponential integrators, ...)
- Concurrent work: **Denoising diffusion sampler** using PIS objective with Ornstein-Uhlenbeck process X and $X_0 \sim \mathcal{N}(0, \nu^2 I)$.

Q. Zhang and Y. Chen. Fast sampling of diffusion models with exponential integrator. arXiv preprint arXiv:2204.13902, 2022

F. Vargas, A. Ovsianas, D. Fernandes, M. Girolami, N. D. Lawrence, and N. Nüsken. Bayesian learning via neural Schrödinger-Föllmer flows. Statistics and Computing, 33(1):1–22, 2023
DIS vs. PIS

Task 1: Compute expectations $\mathbb{E}[f(X_T^u)] \approx \mathbb{E}[f(Y_0)]$ for $f(x) = ||x||^2$ and $f(x) = ||x||_1$ using trained DIS and PIS with varying number of steps in the Euler-Maruyama scheme.



Julius Berner

An optimal control perspective on diffusion-based generative modeling

DIS vs. PIS

Task 2: Compute the log-normalizing constant log \mathcal{Z} for a 2-dim. Gaussian mixture model (GMM) with 9 modes on a 3 × 3 grid, the 10-dimensional Funnel distribution,

$$\rho(x) = \mathcal{N}(x_1; 0, \nu^2) \prod_{i=2}^{10} \mathcal{N}(x_i; 0, e^{x_1}),$$

and the double well (DW) example (using importance sampling in path space).



An optimal control perspective on diffusion-based generative modeling

Summary



Thank you for your attention!

julius.berner@univie.ac.at, lorenz.richter@dida.do, karenu@meta.com