An optimal control perspective on diffusion-based generative modeling

Julius Berner*, Lorenz Richter*, Karen Ullrich

(*equal contribution)

NeurIPS 2022

December, 2022

















SDE-based generative modeling

• Model: Stochastic differential equation (SDE)

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

SDE-based generative modeling

• Model: Stochastic differential equation (SDE)

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

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

SDE-based generative modeling

• Model: Stochastic differential equation (SDE)

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

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



Berner, Richter, Ullrich

An optimal control perspective on diffusion-based generative modeling

PDE for the log-density



PDE for the log-density



PDE for the log-density



Lemma (Hamilton-Jacobi-Bellman equation for log-density)

Define $V(x,t) \coloneqq -\log \bar{p}_X$, where $\bar{p}_X(x,t) \coloneqq p_X(x,T-t)$ denotes the time-reversal. Then

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

Applications



Deriving the ELBO

Theorem (Verification theorem)

Let Y^u be defined by

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

Then it holds that

$$V(Y_0^u, 0) = \min_{u \in \mathcal{U}} \mathbb{E}\left[\underbrace{\int_0^T \left(\operatorname{div}(\mu) + \frac{1}{2} \|u\|^2\right) (Y_s^u, s) \, \mathrm{d}s}_{running \ costs} + \underbrace{V(Y_T^u, T)}_{terminal \ cost} \middle| Y_0^u \right],$$

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

Deriving the ELBO

Theorem (ELBO)

Let Y^u be defined by

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

Then it holds that

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

where the unique minimum is attained by $u^* \coloneqq \sigma^\top \nabla \log p_X = \sigma^\top \nabla \log p_{Y^{u^*}}$.

• **Reparametrization:** $\mu \coloneqq \sigma u - f$ yields the SDEs

 $\mathrm{d} Y_s = f(Y_s,s) \, \mathrm{d} s + \sigma(s) \, \mathrm{d} B_s, \quad Y_0 \sim \mathcal{D}, \qquad \mathrm{d} X^u_s = \left(\bar{\sigma} \, \bar{u} - \bar{f} \right) (X^u_s,s) \, \mathrm{d} s + \bar{\sigma}(s) \, \mathrm{d} B_s.$

• **Reparametrization:** $\mu \coloneqq \sigma u - f$ yields the SDEs

 $\mathrm{d} Y_s = f(Y_s,s) \, \mathrm{d} s + \sigma(s) \, \mathrm{d} B_s, \quad Y_0 \sim \mathcal{D}, \qquad \mathrm{d} X^u_s = \left(\bar{\sigma} \, \bar{u} - \bar{f} \right) (X^u_s,s) \, \mathrm{d} s + \bar{\sigma}(s) \, \mathrm{d} B_s.$



• **Reparametrization:** $\mu \coloneqq \sigma u - f$ yields the SDEs

 $\mathrm{d} Y_s = f(Y_s,s) \, \mathrm{d} s + \sigma(s) \, \mathrm{d} B_s, \quad Y_0 \sim \mathcal{D}, \qquad \mathrm{d} X^u_s = \left(\bar{\sigma} \, \bar{u} - \bar{f} \right) (X^u_s,s) \, \mathrm{d} s + \bar{\sigma}(s) \, \mathrm{d} B_s.$

• **Time-reversal:** If $u = u^* = \sigma^\top \nabla \log p_Y$, one can achieve $X_T^u \sim \mathcal{D}$ by setting $X_0^u \sim Y_T$.



• **Reparametrization:** $\mu \coloneqq \sigma u - f$ yields the SDEs

 $\mathrm{d} Y_s = f(Y_s,s) \, \mathrm{d} s + \sigma(s) \, \mathrm{d} B_s, \quad Y_0 \sim \mathcal{D}, \qquad \mathrm{d} X^u_s = \left(\bar{\sigma} \, \bar{u} - \bar{f} \right) (X^u_s,s) \, \mathrm{d} s + \bar{\sigma}(s) \, \mathrm{d} B_s.$

• **Time-reversal:** If $u = u^* = \sigma^\top \nabla \log p_Y$, one can achieve $X_T^u \sim \mathcal{D}$ by setting $X_0^u \sim Y_T$.



• Denoising score matching objective: The verification theorem implies

$$-\log p_{X_{\mathcal{T}}^{u}}(Y_{0}) \leq \frac{T}{2} \mathbb{E} \left[\left\| u(Y_{\tau}, \tau) - \sigma^{\top}(\tau) \nabla \log p_{Y_{\tau}|Y_{0}}(Y_{\tau}|Y_{0}) \right\|^{2} \right] + \text{const.},$$

Applications



Sampling from densities

• Application: Data distribution \mathcal{D} admits the density $\frac{\rho}{\mathcal{Z}}$ where ρ is known. We cannot use denoising score matching since we do not have access to samples.

Sampling from densities

Application: Data distribution D admits the density ρ/Z where ρ is known. We cannot use denoising score matching since we do not have access to samples.

Corollary (Sampling from densities)

The verification theorem implies

$$-\log \mathcal{Z} = \min_{\overline{u} \in \mathcal{U}} \mathbb{E}\left[\left. \int_0^T \left(\operatorname{div}(\overline{f}) + \frac{1}{2} \|\overline{u}\|^2 \right) (X^u_s, s) \, \mathrm{d}s + \log p_{Y_T}(X^u_0) - \log \rho(X^u_T) \right| Y_0 \right],$$

where the minimum is attained by $u^* = \sigma^\top \nabla \log p_Y$.

Sampling from densities

Application: Data distribution D admits the density ρ/Z where ρ is known. We cannot use denoising score matching since we do not have access to samples.

Corollary (Sampling from densities)

The verification theorem implies

$$-\log \mathcal{Z} = \min_{\overline{u} \in \mathcal{U}} \mathbb{E}\left[\left. \int_0^T \left(\operatorname{div}(\overline{f}) + \frac{1}{2} \| \overline{u} \|^2 \right) (X^u_s, s) \, \mathrm{d}s + \log p_{Y_T}(X^u_0) - \log \rho(X^u_T) \right| Y_0 \right],$$

where the minimum is attained by $u^* = \sigma^\top \nabla \log p_Y$.

• **Related work:** Path Integral Sampler by Zhang and Chen (2021). Our diffusion-based approach is more flexible (SDE type, initial distribution, probability flow ODE) and stable.

Berner, Richter, Ullrich

Numerical example

• Successful application to tasks in molecular dynamics, e.g.,

$$\rho(x) \coloneqq \exp\left(-\sum_{i=1}^{5} (x_i^2 - 2)^2 - \frac{1}{2} \sum_{i=6}^{50} x_i^2\right)$$

in 50 dimensions with 32 modes.

Numerical example

• Successful application to tasks in molecular dynamics, e.g.,

$$\rho(x) \coloneqq \exp\left(-\sum_{i=1}^{5} (x_i^2 - 2)^2 - \frac{1}{2} \sum_{i=6}^{50} x_i^2\right)$$

in 50 dimensions with 32 modes.



Applications



Thank you for your attention!

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

Berner, Richter, Ullrich

An optimal control perspective on diffusion-based generative modeling