

Diffusion Models and Score Matching

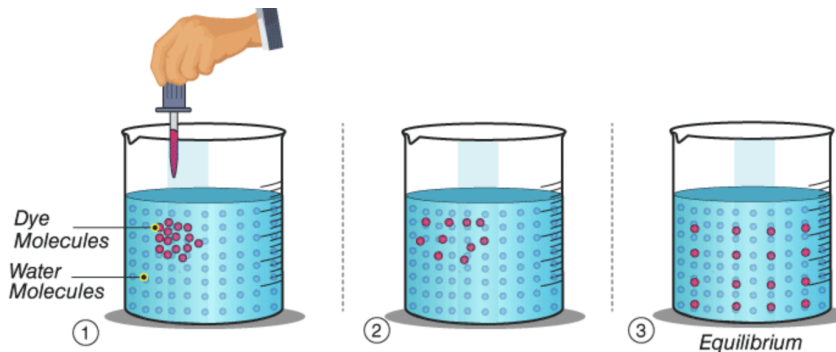
Instructor: Lei Wu ¹

Mathematical Introduction to Machine Learning

Peking University, Fall 2023

¹School of Mathematical Sciences; [Center for Machine Learning Research](#)

What is Diffusion?



Dye molecules **diffuse** throughout the entire space by colliding with water molecules.

Mathematical Modeling: Brownian Motion

- Let $\{x_k\}_{k \geq 0}$ be the trajectory of dye molecules. We can model its dynamics as follows

$$x_{k+1} = x_k + \sqrt{\eta} \xi_k, \quad 0 \leq k \leq N - 1,$$

where $\xi_k \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and η is a small factor ².

² η depends on the temperature, time unit, etc.

Mathematical Modeling: Brownian Motion

- Let $\{x_k\}_{k \geq 0}$ be the trajectory of dye molecules. We can model its dynamics as follows

$$x_{k+1} = x_k + \sqrt{\eta} \xi_k, \quad 0 \leq k \leq N-1,$$

where $\xi_k \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and η is a small factor ².

- Thus, we have after N steps

$$x_{N\eta} = x_0 + \sqrt{\eta} \sum_{k=0}^{N-1} \xi_k \sim \mathcal{N}(x_0, \eta N).$$

² η depends on the temperature, time unit, etc.

Mathematical Modeling: Brownian Motion

- Let $\{x_k\}_{k \geq 0}$ be the trajectory of dye molecules. We can model its dynamics as follows

$$x_{k+1} = x_k + \sqrt{\eta} \xi_k, \quad 0 \leq k \leq N-1,$$

where $\xi_k \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and η is a small factor ².

- Thus, we have after N steps

$$x_{N\eta} = x_0 + \sqrt{\eta} \sum_{k=0}^{N-1} \xi_k \sim \mathcal{N}(x_0, \eta N).$$

- Consider the continuous-time limit: $\eta \rightarrow 0$. Let $t = N\eta$. Then, we have

$$x_{N\eta} \rightarrow X_t \sim \mathcal{N}(X_0, t).$$

² η depends on the temperature, time unit, etc.

Mathematical Modeling: Brownian Motion

- Let $\{x_k\}_{k \geq 0}$ be the trajectory of dye molecules. We can model its dynamics as follows

$$x_{k+1} = x_k + \sqrt{\eta} \xi_k, \quad 0 \leq k \leq N-1,$$

where $\xi_k \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ and η is a small factor ².

- Thus, we have after N steps

$$x_{N\eta} = x_0 + \sqrt{\eta} \sum_{k=0}^{N-1} \xi_k \sim \mathcal{N}(x_0, \eta N).$$

- Consider the continuous-time limit: $\eta \rightarrow 0$. Let $t = N\eta$. Then, we have

$$x_{N\eta} \rightarrow X_t \sim \mathcal{N}(X_0, t).$$

- We call $B_t := X_t - X_0$ Brownian motion.

² η depends on the temperature, time unit, etc.

Important properties of Brownian motion

- After time t , dye molecules only move $O(\sqrt{t})$

$$\mathbb{E}[B_t] = 0, \quad \mathbb{E}[B_t^2] = t.$$

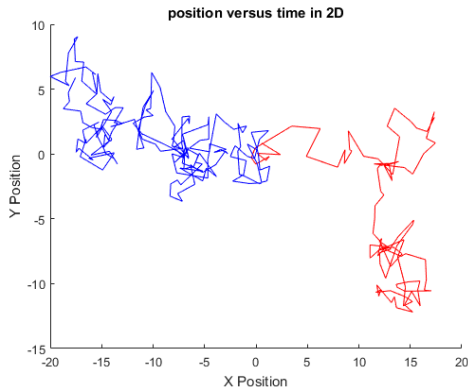


Figure 1: A animation of Brownian motion: https://physics.bu.edu/~duffy/HTML5/brownian_motion.html

Important properties of Brownian motion

- After time t , dye molecules only move $O(\sqrt{t})$

$$\mathbb{E}[B_t] = 0, \quad \mathbb{E}[B_t^2] = t.$$

- $B_t - B_s$ and B_s are independent.

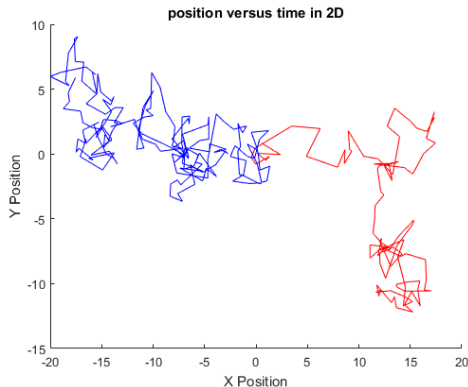


Figure 1: A animation of Brownian motion: https://physics.bu.edu/~duffy/HTML5/brownian_motion.html

Important properties of Brownian motion

- After time t , dye molecules only move $O(\sqrt{t})$

$$\mathbb{E}[B_t] = 0, \quad \mathbb{E}[B_t^2] = t.$$

- $B_t - B_s$ and B_s are independent.
- The trajectory is continuous but **non-differentiable almost everywhere**.

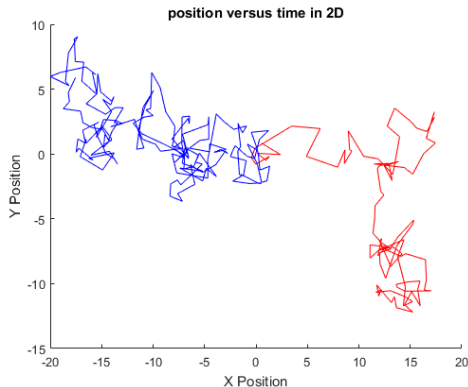


Figure 1: A animation of Brownian motion: https://physics.bu.edu/~duffy/HTML5/brownian_motion.html

General diffusion process (modeled by Ito-SDE)

Consider dye molecules in a force field $f(x, t)$ and the collision is heterogeneous:

- From t to $t + \eta$, the dye molecule moves according to

$$x_{t+\eta} - x_t = \underbrace{f(x_t, t)\eta}_{\text{drift}} + \underbrace{\sigma(x_t, t)\sqrt{\eta}\xi_t}_{\text{diffusion}}.$$

General diffusion process (modeled by Ito-SDE)

Consider dye molecules in a force field $f(x, t)$ and the collision is heterogeneous:

- From t to $t + \eta$, the dye molecule moves according to

$$x_{t+\eta} - x_t = \underbrace{f(x_t, t)\eta}_{\text{drift}} + \underbrace{\sigma(x_t, t)\sqrt{\eta}\xi_t}_{\text{diffusion}}.$$

- Taking $\eta \rightarrow 0$ gives a stochastic differential equation (SDE):

$$dx_t = f(x_t, t) dt + \sigma(x_t, t) dB_t$$

General diffusion process (modeled by Ito-SDE)

Consider dye molecules in a force field $f(x, t)$ and the collision is heterogeneous:

- From t to $t + \eta$, the dye molecule moves according to

$$x_{t+\eta} - x_t = \underbrace{f(x_t, t)\eta}_{\text{drift}} + \underbrace{\sigma(x_t, t)\sqrt{\eta}\xi_t}_{\text{diffusion}}.$$

- Taking $\eta \rightarrow 0$ gives a stochastic differential equation (SDE):

$$dx_t = f(x_t, t) dt + \sigma(x_t, t) dB_t$$

In physics, it is often written (by let $\omega_t = \dot{B}_t$) as

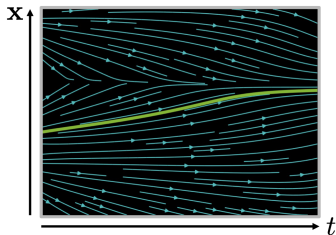
$$\dot{x}_t = f(x_t, t) + \sigma(x_t, t)\omega_t,$$

where ω_t is often referred to as white noise.

A comparison between SDE and ODE ³

Ordinary Differential Equation (ODE):

$$\frac{dx}{dt} = f(x, t) \text{ or } dx = f(x, t)dt$$



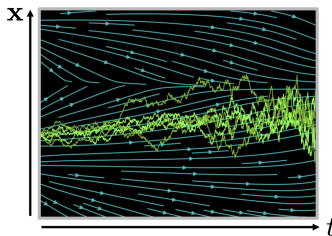
Analytical Solution: $x(t) = x(0) + \int_0^t f(x, \tau) d\tau$

Iterative Numerical Solution: $x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t$

Stochastic Differential Equation (SDE):

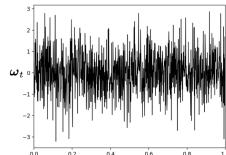
$$\frac{dx}{dt} = \underbrace{f(x, t)}_{\text{drift coefficient}} + \underbrace{\sigma(x, t)\omega_t}_{\text{diffusion coefficient}}$$

$$\left(dx = f(x, t)dt + \sigma(x, t)d\omega_t \right)$$



$$x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t + \sigma(x(t), t)\sqrt{\Delta t}\mathcal{N}(0, \mathbf{I})$$

Wiener Process
(Gaussian
White Noise)



³taken from <https://cvpr2022-tutorial-diffusion-models.github.io/>

Langevin Dynamics

- (Over-damped) **Langevin dynamics** is a special SDE with the drift term given by a potential force $f(x) = -\nabla U(x)$:

$$dx_t = -\nabla U(x_t) dt + \sqrt{2\beta^{-1}} dB_t. \quad (1)$$

Denote by $p_t = p(t, \cdot) = \text{Law}(X_t)$. Then, we have

$$p(t, x) \rightarrow \frac{e^{-\beta U(x)}}{Z} \text{ as } t \rightarrow \infty. \quad (2)$$

Langevin Dynamics

- (Over-damped) **Langevin dynamics** is a special SDE with the drift term given by a potential force $f(x) = -\nabla U(x)$:

$$dx_t = -\nabla U(x_t) dt + \sqrt{2\beta^{-1}} dB_t. \quad (1)$$

Denote by $p_t = p(t, \cdot) = \text{Law}(X_t)$. Then, we have

$$p(t, x) \rightarrow \frac{e^{-\beta U(x)}}{Z} \text{ as } t \rightarrow \infty. \quad (2)$$

- To simulate (1), we can apply the Euler-Maruyama scheme:

$$X_{k+1} = X_k - \nabla U(X_k)\eta + \sqrt{2\beta^{-1}\eta}\xi_k \text{ with } \xi_k \sim \mathcal{N}(0, I_d). \quad (3)$$

Langevin Dynamics

- (Over-damped) **Langevin dynamics** is a special SDE with the drift term given by a potential force $f(x) = -\nabla U(x)$:

$$dx_t = -\nabla U(x_t) dt + \sqrt{2\beta^{-1}} dB_t. \quad (1)$$

Denote by $p_t = p(t, \cdot) = \text{Law}(X_t)$. Then, we have

$$p(t, x) \rightarrow \frac{e^{-\beta U(x)}}{Z} \text{ as } t \rightarrow \infty. \quad (2)$$

- To simulate (1), we can apply the Euler-Maruyama scheme:

$$X_{k+1} = X_k - \nabla U(X_k)\eta + \sqrt{2\beta^{-1}\eta}\xi_k \text{ with } \xi_k \sim \mathcal{N}(0, I_d). \quad (3)$$

- **Ornstein–Uhlenbeck** (OU) process is a simplest SDE given by

$$dx_t = -\theta x_t dt + \sigma dB_t,$$

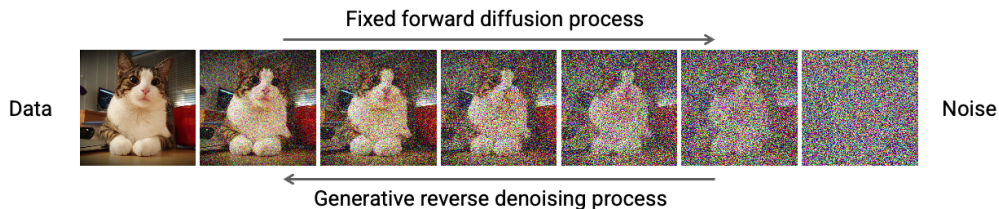
for which $U(x) = \theta\|x\|^2/2$, $\beta^{-1} = \sigma^2/2$. The equilibrium distribution is Gaussian:

$$p_\infty(x) \propto \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right)$$

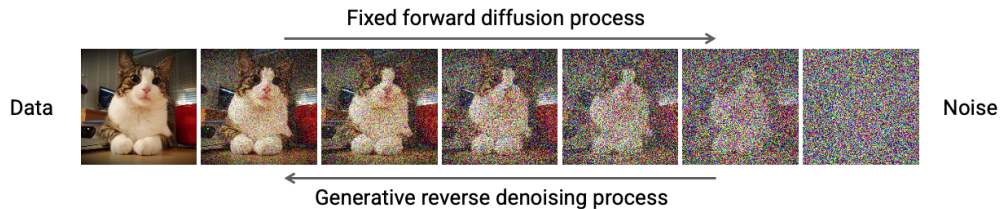
Diffusion Models

In diffusion models

- We first gradually inject noise to a sample until it becomes pure noise. **This is a diffusion process!!**
- The generative models are (probabilistic) inverse of the forward process.



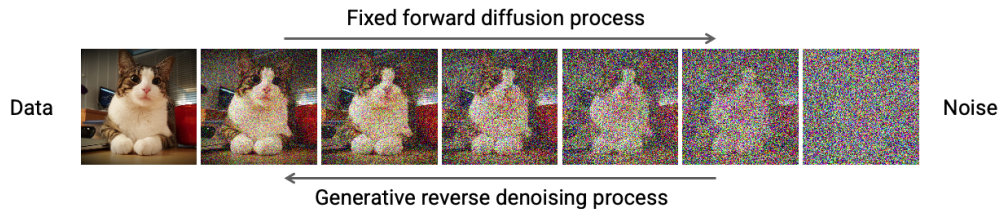
Diffusion Models



Why are diffusion models powerful?

- Guide the learning of reverse **generative** denoise process with the information of a **fixed forward diffusion process**!
- GAN, Normalizing flow, and Variational Autoencoder do not have forward-process information to guide the learning. **[Explain it!]**

Diffusion Models



There are two key issues in diffusion models:

- Construct forward diffusion process.
- Utilize forward information for learning the reverse process.

Denoising Diffusion Probabilistic Models (DDPM) ⁴

- DDPM chooses the following **variance-preserving** forward diffusion process:

$$x_{k+1} = \sqrt{1 - \beta_k} x_k + \sqrt{\beta_k} \xi_k, 0 \leq k \leq N - 1,$$

where $\xi_k \stackrel{iid}{\sim} \mathcal{N}(0, I_d)$.

⁴Jonathan Ho, Ajay Jain, Pieter Abbeel, *Denoising Diffusion Probabilistic Models*, NeurIPS 2020.

Denoising Diffusion Probabilistic Models (DDPM) ⁴

- DDPM chooses the following **variance-preserving** forward diffusion process:

$$x_{k+1} = \sqrt{1 - \beta_k} x_k + \sqrt{\beta_k} \xi_k, 0 \leq k \leq N - 1,$$

where $\xi_k \stackrel{iid}{\sim} \mathcal{N}(0, I_d)$.

- Consider $\beta_t = \beta = o(1)$. Then, we have

$$x_{k+1} = x_k - \frac{\beta x_k}{2} + \sqrt{\beta} \xi_k + o(\beta) \tag{4}$$

When $\beta \rightarrow 0$, we have the forward process is given by an OU process

$$dx_t = -\frac{x_t}{2} dt + dB_t.$$

⁴Jonathan Ho, Ajay Jain, Pieter Abbeel, *Denoising Diffusion Probabilistic Models*, NeurIPS 2020.

Properties of the forward process

- First, the conditional distribution is always Gaussian

$$P_t := x_t | x_0 \sim \mathcal{N}(e^{-t/2}x_0, (1 - e^{-t})I_d) = \mathcal{N}\left(\alpha_t x_0, \sqrt{1 - \alpha_t^2} I_d\right), \quad (5)$$

where $\alpha_t = e^{-t/2}$. We also denote $\sigma_t^2 := 1 - \alpha_t^2$. (Derivation is given on the blackboard.)

Properties of the forward process

- First, the conditional distribution is always Gaussian

$$P_t := x_t|x_0 \sim \mathcal{N}(e^{-t/2}x_0, (1 - e^{-t})I_d) = \mathcal{N}\left(\alpha_t x_0, \sqrt{1 - \alpha_t^2} I_d\right), \quad (5)$$

where $\alpha_t = e^{-t/2}$. We also denote $\sigma_t^2 := 1 - \alpha_t^2$. (Derivation is given on the blackboard.)

- The distribution of x_t can be viewed as the convolution of $P(x_0)$ with a Gaussian smoothing kernel:

$$P_t(x) = \int P_t(x|x_0)P(x_0) \mathrm{d}x_0 = \int P(x_0) \frac{1}{C_t} e^{-\frac{\|x - \alpha_t x_0\|^2}{2(1 - \alpha_t^2)}} \mathrm{d}x_0,$$

where C_t is the normalizing constant.

Properties of the forward process

- First, the conditional distribution is always Gaussian

$$P_t := x_t | x_0 \sim \mathcal{N}(e^{-t/2}x_0, (1 - e^{-t})I_d) = \mathcal{N}\left(\alpha_t x_0, \sqrt{1 - \alpha_t^2} I_d\right), \quad (5)$$

where $\alpha_t = e^{-t/2}$. We also denote $\sigma_t^2 := 1 - \alpha_t^2$. (Derivation is given on the blackboard.)

- The distribution of x_t can be viewed as the convolution of $P(x_0)$ with a Gaussian smoothing kernel:

$$P_t(x) = \int P_t(x|x_0)P(x_0) \mathrm{d}x_0 = \int P(x_0) \frac{1}{C_t} e^{-\frac{\|x - \alpha_t x_0\|^2}{2(1 - \alpha_t^2)}} \mathrm{d}x_0,$$

where C_t is the normalizing constant.

- The forward process converges exponentially fast:

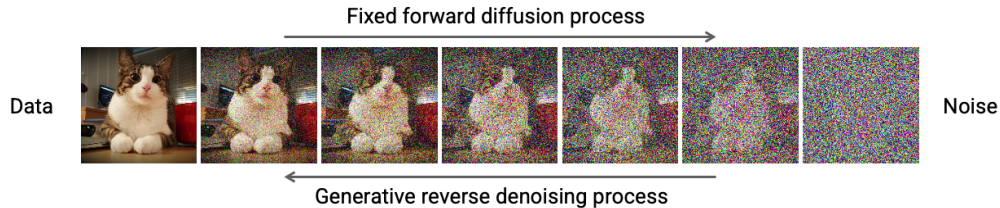
$$D_{\text{KL}}(P_t || \mathcal{N}(0, I_d)) \leq C e^{-t} D_{\text{KL}}(P_0 || \mathcal{N}(0, I_d)),$$

This means we can take a moderately large T such that

$$\text{Law}(x_T) \approx \mathcal{N}(0, I_d).$$

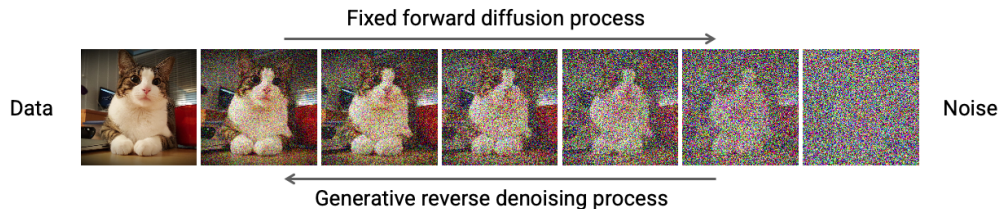
Reversing a Diffusion Process

What do we mean by reversing a diffusion process?



Reversing a Diffusion Process

What do we mean by reversing a diffusion process?



Definition 1

Given a forward process $\{X_t\}_{t \in [0, T]}$, the backward process $\{\tilde{X}_t\}_{t \in [T, 0]}$ is said to be a reverse process of $\{X_t\}_{t \in [0, T]}$ iff

$$\text{Law}(X_t) = \text{Law}(\tilde{X}_{T-t}).$$

Remark: The reverse process may be non-unique.

An Explicit Construction of Reverse Processes

- Consider a large family of diffusion process given by the forward SDE:

$$dx_t = f(x, t) dt + g(t) dB_t, \quad 0 \leq t \leq T.$$

⁵Brian Anderson, *Reverse-time diffusion equation models*, Stochastic Processes and their Applications 12 (1982)

An Explicit Construction of Reverse Processes

- Consider a large family of diffusion process given by the forward SDE:

$$dx_t = f(x, t) dt + g(t) dB_t, \quad 0 \leq t \leq T.$$

- Anderson (1982)**⁵ provided an explicit construction of the reverse SDE:

$$d\tilde{x}_t = [f(\tilde{x}_t, t) - g^2(t) \nabla_x \log p(\tilde{x}, t)] dt + g(t) d\bar{B}_t, \quad t \in [T, 0]$$

where \bar{B}_t is a backward Brownian motion and the time in the above equation is negative. (The proof can be easily completed by checking the Fokker-Planck equation (omitted). We refer to Anderson (1982) for the derivation.)

Total citations [Cited by 390](#)



Scholar articles [Reverse-time diffusion equation models](#)

BDO Anderson - Stochastic Processes and their Applications, 1982

[Cited by 390](#) [Related articles](#) [All 6 versions](#)

⁵Brian Anderson, *Reverse-time diffusion equation models*, Stochastic Processes and their Applications 12 (1982)

Score Matching

- The key quantity for the reverse SDE is the (time-dependent) **score function**

$$\nabla_x \log p(\cdot, \cdot) : \mathbb{R}^d \times [0, T] \mapsto \mathbb{R}^d,$$

Therefore, **we can easily reverse the process if the score function is known.**

Score Matching

- The key quantity for the reverse SDE is the (time-dependent) **score function**

$$\nabla_x \log p(\cdot, \cdot) : \mathbb{R}^d \times [0, T] \mapsto \mathbb{R}^d,$$

Therefore, **we can easily reverse the process if the score function is known.**

- **Model:** Let $s_\theta : \mathbb{R}^d \times [0, T] \mapsto \mathbb{R}^d$ be a neural network to model the score function.

Score Matching

- The key quantity for the reverse SDE is the (time-dependent) **score function**

$$\nabla_x \log p(\cdot, \cdot) : \mathbb{R}^d \times [0, T] \mapsto \mathbb{R}^d,$$

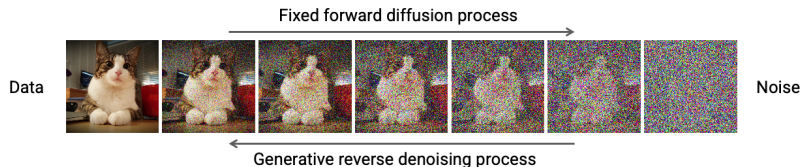
Therefore, **we can easily reverse the process if the score function is known.**

- Model:** Let $s_\theta : \mathbb{R}^d \times [0, T] \mapsto \mathbb{R}^d$ be a neural network to model the score function.
- Training objective:** Let $p_t = p(\cdot, t)$ and

$$L_t(\theta) = \mathbb{E}_{x \sim p_t} [\|s_\theta(x, t) - \nabla_x \log p(x, t)\|^2].$$

Let π be a (weighted) distribution supported on $[0, T]$. Consider the learning via

$$\min_{\theta} L(\theta) := \mathbb{E}_{t \sim \pi} [L_t(\theta)] \quad \textbf{(score matching)}. \quad (6)$$



Score Matching (Cont'd)

- **Why is this objective informative for training?** The problem nearly becomes a sequential of supervised learning: Score matching at different times.
- **Bad News:** $\nabla_x \log p(\cdot, t)$ is unknown. What available are noisy samples $\{x_t\}_{t \in [0, T]}$ generated by the forward process.

Denoising Score Matching

Simplifying the objective using the **log-derivative trick**:

$$\begin{aligned} L_t(\theta) &= \mathbb{E}_{x \sim p_t} [\|s_\theta(x, t) - \nabla_x \log p(x, t)\|^2] \\ &= \mathbb{E}_{x \sim p_t} \|s_\theta(x, t)\|^2 + \mathbb{E}_{x \sim p_t} \|\nabla_x \log p(x, t)\|^2 - 2 \mathbb{E}_{x \sim p_t} \langle s_\theta(x, t), \nabla_x \log p(x, t) \rangle \\ &= \mathbb{E}_{x \sim p_t} \|s_\theta(x, t)\|^2 + \mathbb{E}_{x \sim p_t} \|\nabla_x \log p(x, t)\|^2 - 2 \int_{\mathbb{R}^d} \langle s_\theta(x, t), \nabla_x p(x, t) \rangle dx \\ &= \mathbb{E}_{x \sim p_t} \|s_\theta(x, t)\|^2 + \mathbb{E}_{x \sim p_t} \|\nabla_x \log p(x, t)\|^2 + 2 \int_{\mathbb{R}^d} [\nabla \cdot s_\theta(x, t)] p(x, t) dx. \end{aligned}$$

Note that the conditional distribution $p_t(x|z)$ is tractable.

Denoising Score Matching

Simplifying the objective using the **log-derivative trick**:

$$\begin{aligned} L_t(\theta) &= \mathbb{E}_{x \sim p_t} [\|s_\theta(x, t) - \nabla_x \log p(x, t)\|^2] \\ &= \mathbb{E}_{x \sim p_t} \|s_\theta(x, t)\|^2 + \mathbb{E}_{x \sim p_t} \|\nabla_x \log p(x, t)\|^2 - 2\mathbb{E}_{x \sim p_t} \langle s_\theta(x, t), \nabla_x \log p(x, t) \rangle \\ &= \mathbb{E}_{x \sim p_t} \|s_\theta(x, t)\|^2 + \mathbb{E}_{x \sim p_t} \|\nabla_x \log p(x, t)\|^2 - 2 \int_{\mathbb{R}^d} \langle s_\theta(x, t), \nabla_x p(x, t) \rangle dx \\ &= \mathbb{E}_{x \sim p_t} \|s_\theta(x, t)\|^2 + \mathbb{E}_{x \sim p_t} \|\nabla_x \log p(x, t)\|^2 + 2 \int_{\mathbb{R}^d} [\nabla \cdot s_\theta(x, t)] p(x, t) dx. \end{aligned}$$

Note that the conditional distribution $p_t(x|z)$ is tractable. Noting that $p(x, t) = \int p_t(x|z) p_0(z) dz$, we have

$$\begin{aligned} \int_{\mathbb{R}^d} [\nabla \cdot s_\theta(x, t)] p(x, t) dx &= \int_{\mathbb{R}} p_0(z) dz \int_{\mathbb{R}^d} [\nabla \cdot s_\theta(x, t)] p_t(x|z) dx \\ &= - \int_{\mathbb{R}} p_0(z) dz \int_{\mathbb{R}^d} \langle s_\theta(x, t), \nabla_x p_t(x|z) \rangle dx \\ &= - \int_{\mathbb{R}} p_0(z) dz \int_{\mathbb{R}^d} \langle s_\theta(x, t), \nabla_x \log p_t(x|z) \rangle p_t(x|z) dx \\ &= -\mathbb{E}_{z \sim p_0} \mathbb{E}_{x \sim p_t(\cdot|z)} [\langle s_\theta(x, t), \nabla_x \log p_t(x|z) \rangle] \end{aligned}$$

Denoising Score Matching (Cont'd)

According to the preceding derivation, we have

$$\begin{aligned} L_t(\theta) &= \mathbb{E}_{x \sim p_t} [\|s_\theta(x, t) - \nabla_x \log p(x, t)\|^2] \\ &= \mathbb{E}_{z \sim p_0} \mathbb{E}_{x \sim p_t(\cdot|z)} \|s_\theta(x, t)\|^2 - 2 \mathbb{E}_{z \sim p_0} \mathbb{E}_{x \sim p_t(\cdot|z)} [\langle s_\theta(x, t), \nabla_x \log p_t(x|z) \rangle] + C \\ &= \mathbb{E}_{z \sim p_0} \mathbb{E}_{x \sim p_t(\cdot|z)} [\|s_\theta(x, t) - \nabla_x \log p_t(x|z)\|^2] + C \end{aligned}$$

Consider the DDPM-type forward process and let $\alpha_t = e^{-t/2}$ and $\sigma_t^2 = 1 - e^{-t}$. Then,

$$p_t(x|x_0) \propto \exp\left(-\frac{\|x - \alpha_t x_0\|^2}{2\sigma_t^2}\right).$$

Thus, we have

$$L_t(\theta) = \mathbb{E}_{x_0} \mathbb{E}_{x_t|x_0} \left[\left\| s_\theta(x_t, t) - \frac{x_t - \alpha_t x_0}{\sigma_t^2} \right\|^2 \right]$$

The Denosing Interpretation

In a summary, the total objective becomes

$$L(\theta) = \mathbb{E}_t \mathbb{E}_{x_0} \mathbb{E}_{x_t|x_0} \left[\left\| s_\theta(x_t, t) - \frac{x_t - \alpha_t x_0}{\sigma_t^2} \right\|^2 \right] \quad (7)$$

The Denoising Interpretation

In a summary, the total objective becomes

$$L(\theta) = \mathbb{E}_t \mathbb{E}_{x_0} \mathbb{E}_{x_t|x_0} \left[\left\| s_\theta(x_t, t) - \frac{x_t - \alpha_t x_0}{\sigma_t^2} \right\|^2 \right] \quad (7)$$

Noting that $x_t|x_0 \sim \mathcal{N}(\alpha_t x_0, (1 - \alpha_t)I_d)$, we can rewrite

$$x_t = \alpha_t x_0 + \sqrt{1 - \alpha_t^2} \xi_t = \alpha_t x_0 + \sigma_t \xi_t \quad \text{with } \xi_t \sim \mathcal{N}(0, I_d).$$

Plugging it back into (7) gives the **denoising score matching** objective:

$$L(\theta) = \mathbb{E}_t \mathbb{E}_{x_0} \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)} \left[\left\| s_\theta(x_t, t) - \frac{\xi_t}{\sigma_t} \right\|^2 \right]$$

Training procedure

$$L(\theta) = \mathbb{E}_t \mathbb{E}_{x_0} \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)} \left[\left\| s_\theta(x_t, t) - \frac{\xi_t}{\sigma_t} \right\|^2 \right]$$

Parameterize s_θ with neural networks. Then, SGD of batch size 1 for each iteration updates as follows:

- Step 1: $t \sim \pi, x_0 \sim p_0, \xi_t \sim \mathcal{N}(0, I_d)$
- Step 2: $\alpha_t = e^{-t/2}, x_t = \alpha_t x_0 + \sqrt{1 - \alpha_t} \xi_t$
- Step 3: $\hat{L}(\theta) = \left\| s_\theta(x_t, t) - \frac{\xi_t}{\sigma_t} \right\|^2$
- Step 4: $\theta_{t+1} = \theta_t - \eta \nabla_\theta \hat{L}(\theta_t)$

The Choice of Time Weighting

How to choose π ?

$$L(\theta) = \mathbb{E}_{t \sim \pi} \mathbb{E}_{x_0} \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)} \left[\left\| s_\theta(x_t, t) - \frac{\xi_t}{\sigma_t} \right\|^2 \right]$$

Key observation: When $t \rightarrow 0$, $\sigma_t = \sqrt{1 - e^{-t}} \rightarrow 0$. Loss heavily amplified when sampling t close to 0. High variance!

The Choice of Time Weighting

How to choose π ?

$$L(\theta) = \mathbb{E}_{t \sim \pi} \mathbb{E}_{x_0} \mathbb{E}_{\xi_t \sim \mathcal{N}(0, I_d)} \left[\left\| s_\theta(x_t, t) - \frac{\xi_t}{\sigma_t} \right\|^2 \right]$$

Key observation: When $t \rightarrow 0$, $\sigma_t = \sqrt{1 - e^{-t}} \rightarrow 0$. Loss heavily amplified when sampling t close to 0. High variance!

- Training with **time cut-off** η :

$$\pi = \text{Unif}([\eta, T]).$$

- Variance reduction via **importance sampling**:

$$\pi(t) \propto \frac{1}{\sigma_t^2}.$$

Probability Flow ODE

The reverse SDE is given by

$$d\tilde{x}_t = [f(\tilde{x}_t, t) - g^2(t) \nabla_x \log p(\tilde{x}, t)] dt + g(t) d\bar{B}_t,$$

(Song et al. 2021) showed that the following probability-flow ODE is also a reverse process

$$d\tilde{x}_t = f(\tilde{x}_t, t) - \frac{1}{2}g^2(t) \nabla_x \log p(\tilde{x}, t) dt.$$

For the DDPM-type forward process, it becomes

$$d\tilde{x}_t = -\frac{1}{2}(x_t + \nabla_x \log p(\tilde{x}, t)) dt$$

The probability-flow ODE can be interpreted as a continuous-time normalizing flow (CNF).

A Schematic Comparison

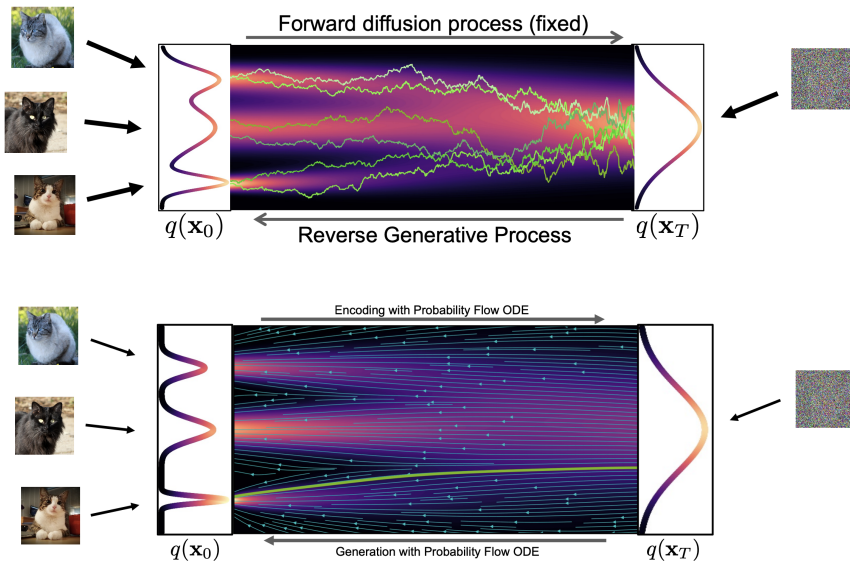


Figure 2: (Up) SDE; (Down) ODE.

Faster Sampling

- For diffusion models, generating new samples needs to discretize the revise-time SDE. It is often very slow as we need to take small step size for reducing discretization error and numerical stability.
 - For SDEs, in general, there does not exist higher-order solver as the trajectory is non-differentiable almost everywhere.
 - Deterministic ODE enables the use of advanced higher-order ODE solvers such as Runge-Kutta, thereby speeding up the generation of new samples.

Exact Likelihood Computation ⁶

- Consider the continuous-time normalizing flow generated by the ODE

$$\dot{x}_t = f(x_t, t), t \in [0, T]$$

with initial condition $x_0 \sim p_0$.

- Consider the flow map $\Phi_T : \mathbb{R}^d \mapsto \mathbb{R}^d$ defined by $\Phi_T(x_0) = x_T$. Then, we have the log-likelihood of p_T satisfies (with the derivation left as homework)

$$\log p_T(x_T) = \log p_0(x_0) - \int_0^T \nabla \cdot f(x_t, t) dt.$$

- In practice, $\nabla \cdot f(x, t)$ is estimated using the Skilling-Hutchinson trace estimator:

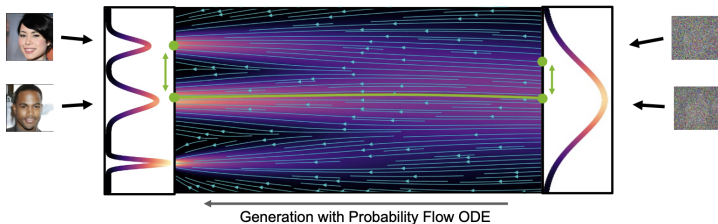
$$\underbrace{\nabla_x \cdot f(x_t, t)}_{\text{need } d \text{ gradients}} = \mathbb{E}_\epsilon[\epsilon^T \nabla_x f(x, t) \epsilon] \approx \frac{1}{m} \sum_{j=1}^m \underbrace{\epsilon_j \cdot \nabla_x (f(x, t) \cdot \epsilon_j)}_{\text{need only } m \text{ gradients}},$$

where $\mathbb{E}[\epsilon] = 0$ and $\mathbb{E}[\epsilon \epsilon^\top] = I_d$.

⁶For generality, we assume normal time direction in this slide.

Manipulating the Latent Space

Interpolation



Controllable Generalization

Generate one-class of samples

class: bird



class: deer



Text to Images

You
请画一幅暴雪中的长城景色

ChatGPT

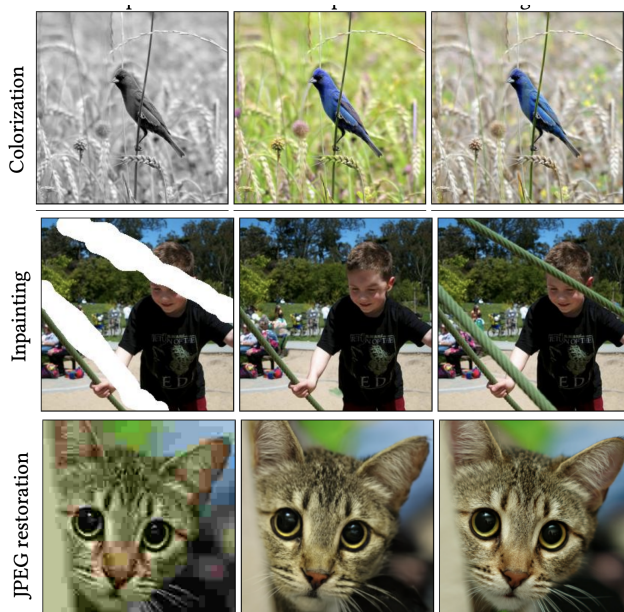


You
请用八大山人的风格重绘这幅画，不要改变内容，只改变风格。请注意留白和落款，不用输出任何文字描述。

ChatGPT



Some Classical Tasks



Controllable Generalization via Diffusion Models

- Controlled generalization can be modeled as sampling from $P(x|y)$ where y denotes the control factor.

Controllable Generalization via Diffusion Models

- Controlled generalization can be modeled as sampling from $P(x|y)$ where y denotes the control factor.
- Baye's rule: $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$. Accordingly, the score function:

$$\nabla \log P(x|y) = \nabla \log P(y|x) + \nabla \log P(x) \quad (8)$$

where the **normalizing constant disappears**.

Controllable Generalization via Diffusion Models

- Controlled generalization can be modeled as sampling from $P(x|y)$ where y denotes the control factor.
- Baye's rule: $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$. Accordingly, the score function:

$$\nabla \log P(x|y) = \nabla \log P(y|x) + \nabla \log P(x) \quad (8)$$

where the **normalizing constant disappears**.

- We can use Langevin dynamics to simulate it. But we can directly couple (8) with the reverse-time SDE or ODE:

$$\begin{aligned} \dot{\tilde{x}}_t &= f(\tilde{x}_t, t) - \frac{1}{2}g(t)^2 \log p_t(\tilde{x}|y) \\ &= f(\tilde{x}_t, t) - \frac{1}{2}\nabla_x [\log p(\tilde{x}_t, t) + \log p(y|\tilde{x}_t)]. \end{aligned}$$

Controllable Generalization via Diffusion Models

- Controlled generalization can be modeled as sampling from $P(x|y)$ where y denotes the control factor.
- Baye's rule: $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$. Accordingly, the score function:

$$\nabla \log P(x|y) = \nabla \log P(y|x) + \nabla \log P(x) \quad (8)$$

where the **normalizing constant disappears**.

- We can use Langevin dynamics to simulate it. But we can directly couple (8) with the reverse-time SDE or ODE:

$$\begin{aligned} \dot{\tilde{x}}_t &= f(\tilde{x}_t, t) - \frac{1}{2}g(t)^2 \log p_t(\tilde{x}|y) \\ &= f(\tilde{x}_t, t) - \frac{1}{2}\nabla_x [\log p(\tilde{x}_t, t) + \log p(y|\tilde{x}_t)]. \end{aligned}$$

- Therefore, as long as we have a good “classifier” $p(y|x)$, then we can couple it with the unconditional model $s_\theta(x, t) \approx \nabla \log p(x, t)$ in a very **easy and principled** approach.

Controllable Generalization via Diffusion Models

- Controlled generalization can be modeled as sampling from $P(x|y)$ where y denotes the control factor.
- Baye's rule: $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$. Accordingly, the score function:

$$\nabla \log P(x|y) = \nabla \log P(y|x) + \nabla \log P(x) \quad (8)$$

where the **normalizing constant disappears**.

- We can use Langevin dynamics to simulate it. But we can directly couple (8) with the reverse-time SDE or ODE:

$$\begin{aligned} \dot{\tilde{x}}_t &= f(\tilde{x}_t, t) - \frac{1}{2}g(t)^2 \log p_t(\tilde{x}|y) \\ &= f(\tilde{x}_t, t) - \frac{1}{2}\nabla_x [\log p(\tilde{x}_t, t) + \log p(y|\tilde{x}_t)]. \end{aligned}$$

- Therefore, as long as we have a good “classifier” $p(y|x)$, then we can couple it with the unconditional model $s_\theta(x, t) \approx \nabla \log p(x, t)$ in a very **easy and principled** approach.
- One unconditional models for all tasks.**

Connection with energy-based models

- In EBM, $p(x) = e^{-U(x)}/Z$. We learn a potential **energy** $V_\theta(x) \approx U(x)$.
- In score-based models, we learn $s_\theta(x) \approx \nabla_x \log p(x) = -\nabla_x U(x)$, i.e., the **force**.
- With the score functions (aka. the force field), we can also recover samples by running Langevin dynamics

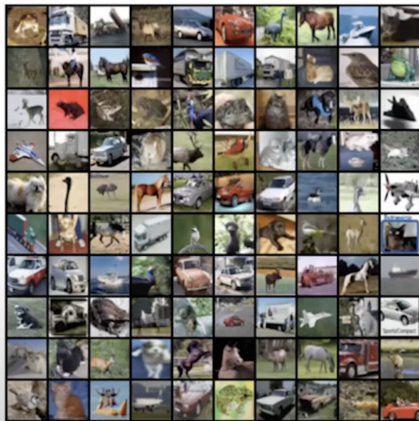
$$dx_t = -\nabla_x \log p(x_t) dt + \sqrt{2} dB_t.$$

But its performance is notorious and consequently, using the reverse-time SDE/ODE is always much better.

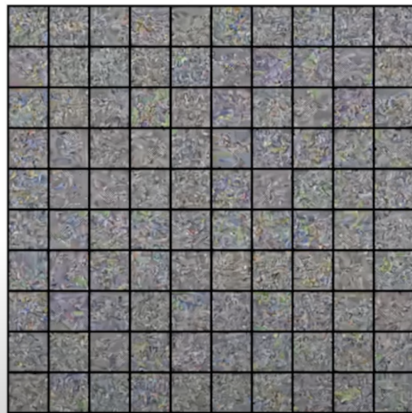
Naive Score Matching + Langevin Dynamics

Often, the learned score function is useless when simulating Langevin dynamics.

CIFAR-10 data



Model samples



Explanation

The learned score function are inaccurate in the low-density region.

$$L(\theta) = \mathbb{E}_x \|s_\theta(x) - \nabla \log p(x)\|^2.$$

