

# Diffusion Model and Stochastic Differential Equation

Chuwei Wang

School of Mathematical Sciences

**Peking University**

2022.10.21

# Outline

- (i) Background of generative model
- (ii) Discrete diffusion model
- (iii) Continuous diffusion model
- (iv) Diffusion model with Bridge
- (v) Summary

# Generative Model

Generative Model plays an important role in numerous tasks, such as computer vision and natural language processing.

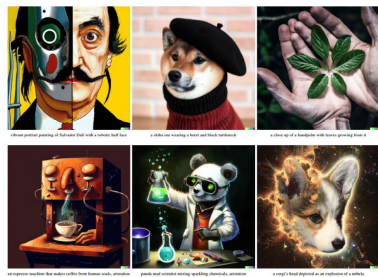


Figure: 1



Figure: 2

## Formulation of generative tasks

- INPUT: Samplings  $\{x_i\}$  from an unknown distribution  $\mu$ .
- OUTPUT: Generated  $\{\hat{x}_i\}$  whose distribution  $\hat{\mu} \approx \mu$ .

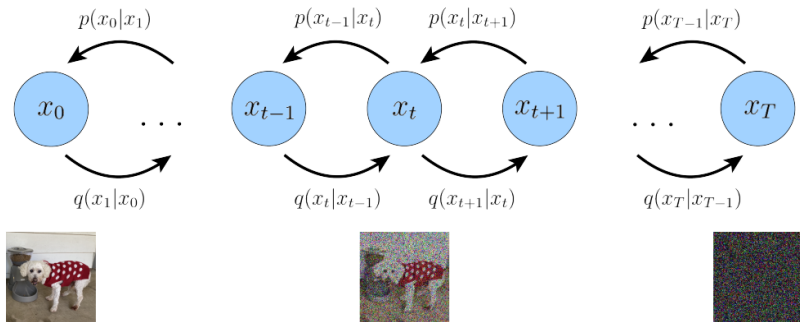
## Methodology

- Based on the samplings: to deal with the empirical distribution  $\frac{1}{n} \sum_i \delta_{x_i}$
- Based on an known distribution  $\nu$  (usually Gaussian): try to transform  $\nu$  into  $\mu$ .

# How to transform $\nu$ into $\mu$ ?

- (straightforward) Look for a mapping  $F_\theta$  such that the  $\mu_{F_\theta(X)} \approx \mu, X \sim \nu$ .
- Find a sequence of mapping  $\{F_\theta^i\}_{i=1}^n$ , such that  $\mu_{F_\theta^n \circ \dots \circ F_\theta^1(X)} \approx \mu, X \sim \nu$ .
- Find a continuous path  $F(t)$  in the space of probability measure, such that  $F(0) = \mu, F(1) = \nu$ .

# Discrete diffusion model



$x_0$  : sampling from the unknown distribution  $\mu (: p(x))$ .

Towards the left: adding noise.  $q(x_{t+1}|x_t) = N(\sqrt{\alpha_t}x_t, (1 - \alpha_t)I)$

Towards the right: denoising.  $p(x_{t-1}|x_t)$  is parameterized and to be learnt.

# Training diffusion model

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \end{aligned}$$

$$\mathbf{x}_t \sim N(\sqrt{\bar{\alpha}_t}, (1 - \bar{\alpha}_t)I), \quad \bar{\alpha}_t := \prod_{i=1}^t \alpha_i$$

Parameterized  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  as  $N(\mu_\theta(\mathbf{x}_t), \sigma(t)^2 I)$ , and further write  $\mu_\theta(\mathbf{x}_t)$  as  $C_1(\alpha)\mathbf{x}_t + C_2(\alpha)\mathbf{x}_\theta(\hat{\mathbf{x}}_t, t)$ . We have

$$\begin{aligned} &\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} \left[ \|\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2 \right] \end{aligned}$$

# Training diffusion model

- (i) Sample  $x_0$  from the latent distribution.
- (ii) Generate the sequence  $\{x_t\}$  by adding gaussian noise.
- (iii) Parameterize  $x_\theta(x_t, t)$ , the estimation of  $x_0$  at the state  $(x_t, t)$ .
- (iv)  $\operatorname{argmin}_\theta \sum_t \mathbb{E}_{q(x_t|x_0)} \lambda(t) \|x_\theta - x_0\|_2^2$ .

## Generating

Assumption:  $p_{x_T} \approx N(0, I)$ .

- (i) Generate  $z \sim N(0, I)$ ,
- (ii) Generate  $\hat{x}$  through the denoising sequence:  
 $p(z)p_\theta(x_{T-1}|x_T)\dots p_\theta(x_1|x_2)p_\theta(x_0|x_1)$ .



# Training diffusion model

## Remark

- (1) There are other parametrizations, eg. the scope ( $\nabla \log p(x_t)$ ) and the source noise, while the loss function remains similar (weighted  $L^2$ ).
- (2) The choice for the hyperparameter  $\{\alpha_t\}$ :
  - Too large:  $\{x_t\}$  converges to  $N(0, I)$  slowly.
  - Too small: Approximating  $p(x_{t-1}|x_t)$  with normal distribution would bring much error.

Published as a conference paper at ICLR 2021

---

## SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

**Yang Song\***

Stanford University

yangsong@cs.stanford.edu

**Jascha Sohl-Dickstein**

Google Brain

jaschasd@google.com

**Diederik P. Kingma**

Google Brain

durk@google.com

**Abhishek Kumar**

Google Brain

abhishk@google.com

**Stefano Ermon**

Stanford University

ermon@cs.stanford.edu

**Ben Poole**

Google Brain

pooleb@google.com

# Basic Knowledge

- [Def] **Ito Integral**:  $\int_0^T f(t)dW_t := \lim_{\Delta \rightarrow 0} \sum f(t_n)(W_{t_{n+1}} - W_{t_n})$ ,  $W_t$  is standard Brownian Motion in  $R^d$ .
- [Def]SDE(Diffusion Process)  $dX_t = b(X_t, t)dt + \sigma(X_t, t)dW_t$ :  
$$X_T - X_0 = \int_{t=0}^T b(X_t, t)dt + \int_{t=0}^T \sigma(X_t, t)dW_t$$
- $\{X_t\}$  is a stochastic process,  $b: R^d$ -valued,  $\sigma: R^{d \times m}$ -valued,  $W_t$ : standard Brownian Motion in  $R^m$ .

# Continuous diffusion model

Original diffusion model is the discretization of a continuous version.

## Example

$$x_{t+1} = \sqrt{\alpha_t}x_t + \sqrt{1 - \alpha_t}w_t, \quad w_t \sim N(0, I) \quad (1)$$

$$\iff x_{t+1} - x_t = (\sqrt{\alpha_t} - 1)x_t + \sqrt{1 - \alpha_t}w_t \quad (2)$$

$$\iff \Delta x_t = (\sqrt{\alpha_t} - 1)x_t \Delta t + \sqrt{1 - \alpha_t} \Delta w_t, \quad \Delta t = 1 \quad (3)$$

$$\rightarrow dX_t = (\sqrt{\alpha(t)} - 1)X_t dt + \sqrt{1 - \alpha(t)}dW_t. \quad (4)$$

Original diffusion model corresponds to Ornstein - Uhlenbeck process, which converges to normal distribution as  $t \rightarrow \infty$ .

# Continuous diffusion model

## Theorem

The reverse of the diffusion process characterized by

$dx = f(x, t)dt + g(t)dw$ ,  $x(0) \sim p_0$  is the solution to

$dx = [f(x, t) - g(t)^2 \nabla \log p_t(x)]dt + g(t)d\tilde{w}$ , where  $p_t$  denotes the marginal distribution of  $x$  at time  $t$ .

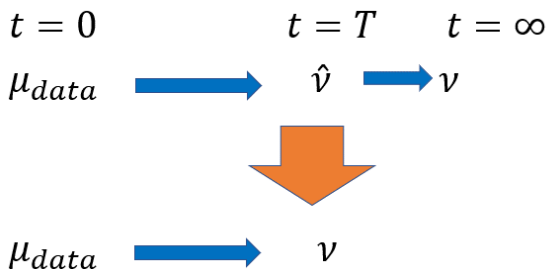
This inspires us to approximate  $\nabla \log p_t(x)$  with an NN function  $s_\theta(x(t), t)$  and reconstruct  $x_0$  from the noise basing on the reverse equation. We could learn  $s_\theta$  through the following optimization.

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t)|\mathbf{x}(0)} \left[ \left\| \mathbf{s}_\theta(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \right\|_2^2 \right] \right\}$$

# Recent Advancement

- (1) Apply computational techniques to improve efficiency and accuracy, e.g. high order scheme and adaptive sampling for solving SDE, predictor-corrector scheme for higher accuracy.
- (2) Some researcher observe that the marginal probability density  $\{p_t(x)\}_{t=0}^T$  shares the same trajectory with a deterministic ODE defined as  $dx = [f(x, t) - \frac{1}{2}g(t)^2 \nabla \log p_t(x)]dt$ . Thus, one can avoid simulating SDEs for generating data.

# Diffusion model with bridge



In previous SDE, the distribution converges to  $N(0, I)$  when  $t \rightarrow \infty$ . Thus, we have to solve the SDE for a large time scale, which is very time-consuming.

---

# Deep Generative Learning via Schrödinger Bridge

---

Gefei Wang<sup>1</sup> Yuling Jiao<sup>2</sup> Qian Xu<sup>3</sup> Yang Wang<sup>1,4</sup> Can Yang<sup>1,4</sup>

## Schrodinger Bridge Problem

Let  $\Omega = C([0, 1], \mathbb{R}^d)$ .  $\mathcal{P}(\Omega)$  is the space of probability measure on the path space  $\Omega$ .  $P_\tau$  is a Brownian motion. Given  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ , find  $Q^* \in \operatorname{argmin}_{\mathcal{P}(\Omega)} \mathbb{D}_{KL}(Q || P_\tau)$  s.t.  $Q_0 = \mu, Q_1 = \nu$ .



# Diffusion model with bridge

Fortunately, Schrodinger Bridge Problem has been well-studied in math. Applying the result about the properties of the solution, this work derives a finite-time generative method.

**Theorem 3** Define the density ratio  $f(\mathbf{x}) = \frac{q_\sigma(\mathbf{x})}{\Phi_{\sqrt{\tau}}(\mathbf{x})}$ . Then for the SDE

$$d\mathbf{x}_t = \tau \nabla \log \mathbb{E}_{\mathbf{z} \sim \Phi_{\sqrt{\tau}}}[f(\mathbf{x}_t + \sqrt{1-t}\mathbf{z})]dt + \sqrt{\tau}d\mathbf{w}_t \quad (4)$$

with initial condition  $\mathbf{x}_0 = \mathbf{0}$ , we have  $\mathbf{x}_1 \sim q_\sigma(\mathbf{x})$ .

And, for the SDE

$$d\mathbf{x}_t = \sigma^2 \nabla \log q_{\sqrt{1-t}\sigma}(\mathbf{x}_t)dt + \sigma d\mathbf{w}_t \quad (5)$$

with initial condition  $\mathbf{x}_0 \sim q_\sigma(\mathbf{x})$ , we have  $\mathbf{x}_1 \sim p_{\text{data}}(\mathbf{x})$ .

$$\Phi_\sigma(\cdot) : N(0, \sigma^2 I), \quad q_\sigma = p_{\text{data}} * \Phi_\sigma.$$

---

# Diffusion-based Molecule Generation with Informative Prior Bridges

---

**Lemeng Wu\***

University of Texas at Austin

lmwu@cs.utexas.edu

**Chengyue Gong\***

University of Texas at Austin

cygong@cs.utexas.edu

**Xingchao Liu**

University of Texas at Austin

xcliu@cs.utexas.edu

**Mao Ye**

University of Texas at Austin

my21@cs.utexas.edu

**Qiang Liu**

University of Texas at Austin

lqiang@cs.utexas.edu

# How to design SDE?

## Target

Transform  $\nu$  into  $\mu$ . (Create an SDE  $dx = f(x, t)dt + g(x, t)dw$  such that  $x_0 \sim \nu$ ,  $x_1 \sim \mu$ .)

- (i) Try to transform  $\nu$  into  $\delta_{x'}$ .
- (ii) Every  $x \in \mathbb{R}^d$  is related to a transformation (a mapping/trajectory from  $\nu$  towards  $\delta_x$ ). Take expectation of these transformations over  $x \sim \mu$  and derive a transformation from  $\nu$  to  $\mu$ .

For (i), get inspiration from the gradient field of Lyapunov function and consider  $f(x, t) = -\alpha_t \nabla U(x) + v(x, t)$ , where  $U$  is a Lyapunov function at  $x'$ ,  $\alpha_t$  controls the step size of gradient flow,  $v(x, t)$  is a perturbation term.

# Summary

- VAE  $\rightarrow$  HVAE and diffusion model
- discrete  $\rightarrow$  continuous
- infinite time scale  $\rightarrow$  finite time convergence

## Inspiration for FermiNet

For Diffusion Monte Carlo, a large proportion of iterations play the role of "tending to infinity", which is very time-consuming. It might be helpful if we could design an SDE dynamic such that  $\psi_I$  evolves into  $\psi^*$  within finite time.

Thanks!