

# Denoising diffusion probabilistic models and normalizing flows

June 18, 2025

# Reminder

- ① **Solve problems:** In total, 8 – 10 questions shall be solved.
- ② **one course project:** Select 1 topic to conduct numerical experiment, and write a short but complete report.
  - length: no requirement. 4-6 pages recommended.
  - problem description and goal
  - method and algorithm
  - details about numerical experiment
- ③ give a 5-10 min presentation.

Remark:

- ① Two persons can work together on the course report. In this case, the project should be more comprehensive. Each should contribute equally to the experiment and to the writing.
- ② General code of conduct for scientific writing applies (e.g. use of materials of this course, online materials, and ChatGPT)!

Part 1: Recall the previous lecture

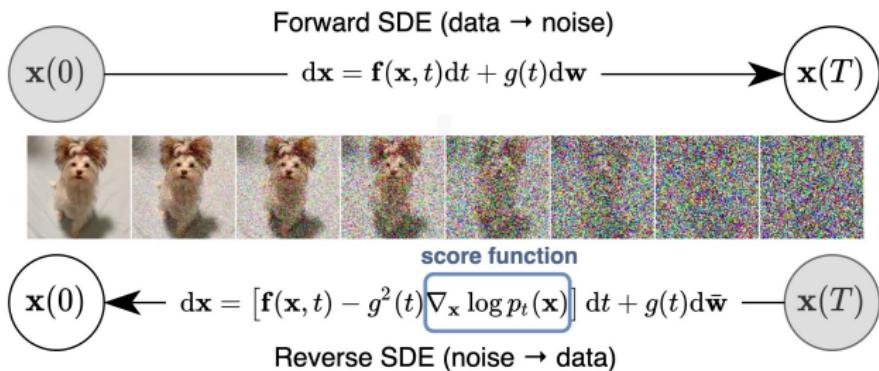
# Generative modeling

Different approaches:

- ① Variational AutoEncoders (VAEs)
- ② Generative Adversarial Networks (GANs)
- ③ Normalizing Flows (NFs)
- ④ Diffusion generative models
- ⑤ Flow-matching generative models

# Score-based diffusion models

- ① choose a forward process  $X_t$   
linear SDE  $\Rightarrow p(\cdot, T)$  is (approximately) Gaussian  $p_{\text{prior}}$ .
- ② learn SDE of the backward process  $Y_t$
- ③ generate new data: sample  $Y_0 \sim p_{\text{prior}}$  and simulate the backward process to get  $Y_T$



# Flow-matching generative models

- ➊ Target density  $p_1 = p_{\text{target}}$  on  $\mathbb{R}^d$ .
- ➋ Prior density  $p_0$  on  $\mathbb{R}^d$ , typically a Gaussian density, is chosen.
- ➌ Define  $p(\cdot, t)$  as the probability density of

$$X_t = (1 - t)X_0 + tX_1, \quad \text{where } X_0 \sim p_0 \text{ and } X_1 \sim p_1. \quad (1)$$

Idea: learn an ODE

$$\frac{dY_t}{dt} = u(Y_t, t), \quad t \in [0, 1] \quad (2)$$

such that, when  $Y_0 \sim p_0$ , then  $Y_t \sim p(\cdot, t)$  for any  $t \in [0, 1]$ .

## Part 2: Denoising diffusion probabilistic models (DDPMs)

---

## Denoising Diffusion Probabilistic Models

---

**Jonathan Ho**

UC Berkeley

[jonathanho@berkeley.edu](mailto:jonathanho@berkeley.edu)

**Ajay Jain**

UC Berkeley

[ajayj@berkeley.edu](mailto:ajayj@berkeley.edu)

**Pieter Abbeel**

UC Berkeley

[pabbeel@cs.berkeley.edu](mailto:pabbeel@cs.berkeley.edu)

### Abstract

We present high quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>

citation statistics: > 23000.

# Basics

- ➊ Gaussian density with mean  $\mu \in \mathbb{R}^d$  and covariance  $\Sigma \in \mathbb{R}^{d \times d}$ :

$$\mathcal{N}(x; \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} (\det \Sigma)^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}, \quad x \in \mathbb{R}^d. \quad (3)$$

# Basics

- ① Gaussian density with mean  $\mu \in \mathbb{R}^d$  and covariance  $\Sigma \in \mathbb{R}^{d \times d}$ :

$$\mathcal{N}(x; \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} (\det \Sigma)^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}, \quad x \in \mathbb{R}^d. \quad (3)$$

- ②  $f : \mathbb{R} \rightarrow \mathbb{R}$  is convex. Jensen's inequality, for any  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

$$f(\mathbb{E}_\pi(g(x))) = f\left(\int_{\mathbb{R}^d} g(x)\pi(x)dx\right) \leq \mathbb{E}_\pi(f(g(x))). \quad (4)$$

# Basics

- ① Gaussian density with mean  $\mu \in \mathbb{R}^d$  and covariance  $\Sigma \in \mathbb{R}^{d \times d}$ :

$$\mathcal{N}(x; \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} (\det \Sigma)^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}, \quad x \in \mathbb{R}^d. \quad (3)$$

- ②  $f : \mathbb{R} \rightarrow \mathbb{R}$  is convex. Jensen's inequality, for any  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

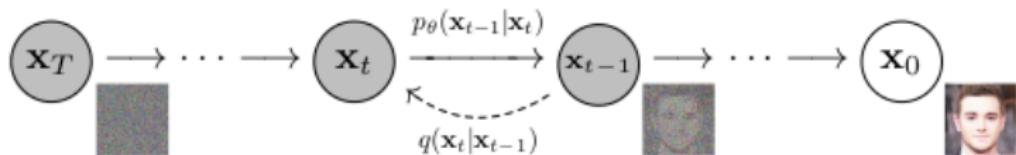
$$f(\mathbb{E}_\pi(g(x))) = f\left(\int_{\mathbb{R}^d} g(x)\pi(x)dx\right) \leq \mathbb{E}_\pi(f(g(x))). \quad (4)$$

- ③ KL divergence

$$D_{KL}(Q_1 | Q_2) := \mathbb{E}_{Q_1}\left(\log \frac{Q_1}{Q_2}\right). \quad (5)$$

- $D_{KL}(Q_1 | Q_2) \geq 0$
- $D_{KL}(Q_1 | Q_2) \neq D_{KL}(Q_2 | Q_1)$

# Markov chains



# Markov chains

- ① Data distribution:  $q_0(x)dx$ .

# Markov chains

- ① Data distribution:  $q_0(x)dx$ .
- ② Forward process:  $x^{(1)}, \dots, x^{(N)} \in \mathbb{R}^d$ , given  $x^{(0)} \in \mathbb{R}^d$ . Joint probability:

$$q(x^{(1:N)} | x^{(0)}) = \prod_{k=1}^N q(x^{(k)} | x^{(k-1)}),$$

where the transition density is

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

# Markov chains

- ① Data distribution:  $q_0(x)dx$ .
- ② Forward process:  $x^{(1)}, \dots, x^{(N)} \in \mathbb{R}^d$ , given  $x^{(0)} \in \mathbb{R}^d$ . Joint probability:

$$q(x^{(1:N)} | x^{(0)}) = \prod_{k=1}^N q(x^{(k)} | x^{(k-1)}),$$

where the transition density is

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

- ③ Reverse process:

$$p_\theta(x^{(0:N)}) = p(x^{(N)}) \prod_{k=1}^N p_\theta(x^{(k-1)} | x^{(k)}),$$

where  $p(x^{(N)})$  is a prior,  $\theta$  is parameter, and

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \Sigma_\theta(x^{(k)}, k)\right).$$

# Markov chains

- ① Data distribution:  $q_0(x)dx$ .
- ② Forward process:  $x^{(1)}, \dots, x^{(N)} \in \mathbb{R}^d$ , given  $x^{(0)} \in \mathbb{R}^d$ . Joint probability:

$$q(x^{(1:N)} | x^{(0)}) = \prod_{k=1}^N q(x^{(k)} | x^{(k-1)}),$$

where the transition density is

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

- ③ Reverse process:

$$p_\theta(x^{(0:N)}) = p(x^{(N)}) \prod_{k=1}^N p_\theta(x^{(k-1)} | x^{(k)}),$$

where  $p(x^{(N)})$  is a prior,  $\theta$  is parameter, and

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \Sigma_\theta(x^{(k)}, k)\right).$$

- ④  $p_\theta(x^{(0)}) = \int p_\theta(x^{(0:N)}) dx^{(1:N)}.$

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\mathbb{E}_{q_0}(-\log p_\theta(x^{(0)}))$$

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\begin{aligned} & \mathbb{E}_{q_0}(-\log p_\theta(x^{(0)})) \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int p_\theta(x^{(0:N)}) dx^{(1:N)}\right)\right] \end{aligned}$$

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\begin{aligned} & \mathbb{E}_{q_0}(-\log p_\theta(x^{(0)})) \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int p_\theta(x^{(0:N)}) dx^{(1:N)}\right)\right] \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})} q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right)\right] \end{aligned}$$

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\begin{aligned} & \mathbb{E}_{q_0}(-\log p_\theta(x^{(0)})) \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int p_\theta(x^{(0:N)}) dx^{(1:N)}\right)\right] \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})} q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right)\right] \\ &\leq \mathbb{E}_{q_0}\left[\int -\log\left(\frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right] \end{aligned}$$

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\begin{aligned} & \mathbb{E}_{q_0}(-\log p_\theta(x^{(0)})) \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int p_\theta(x^{(0:N)}) dx^{(1:N)}\right)\right] \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})} q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right)\right] \\ &\leq \mathbb{E}_{q_0}\left[\int -\log\left(\frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right] \\ &= \mathbb{E}_{\mathbb{Q}}\left(-\log \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) \end{aligned}$$

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\begin{aligned} & \mathbb{E}_{q_0}(-\log p_\theta(x^{(0)})) \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int p_\theta(x^{(0:N)}) dx^{(1:N)}\right)\right] \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})} q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right)\right] \\ &\leq \mathbb{E}_{q_0}\left[\int -\log\left(\frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right] \\ &= \mathbb{E}_{\mathbb{Q}}\left(-\log \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) \\ &= \mathbb{E}_{\mathbb{Q}}\left(-\log p(x^{(N)}) - \sum_{k=1}^N \log \frac{p_\theta(x^{(k-1)} | x^{(k)})}{q(x^{(k)} | x^{(k-1)})}\right) =: L, \end{aligned} \tag{6}$$

where  $\mathbb{E}_{\mathbb{Q}}$  is expectation w.r.t. the density  $q(x^{(0:N)})$  of the forward process.

## Variational bound on negative log-likelihood

Using Jensen's inequality for the convex function  $f(x) = -\log(x)$ ,

$$\begin{aligned} & \mathbb{E}_{q_0}(-\log p_\theta(x^{(0)})) \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int p_\theta(x^{(0:N)}) dx^{(1:N)}\right)\right] \\ &= \mathbb{E}_{q_0}\left[-\log\left(\int \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})} q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right)\right] \\ &\leq \mathbb{E}_{q_0}\left[\int -\log\left(\frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) q(x^{(1:N)} | x^{(0)}) dx^{(1:N)}\right] \\ &= \mathbb{E}_{\mathbb{Q}}\left(-\log \frac{p_\theta(x^{(0:N)})}{q(x^{(1:N)} | x^{(0)})}\right) \\ &= \mathbb{E}_{\mathbb{Q}}\left(-\log p(x^{(N)}) - \sum_{k=1}^N \log \frac{p_\theta(x^{(k-1)} | x^{(k)})}{q(x^{(k)} | x^{(k-1)})}\right) =: L, \end{aligned} \tag{6}$$

where  $\mathbb{E}_{\mathbb{Q}}$  is expectation w.r.t. the density  $q(x^{(0:N)})$  of the forward process.

Idea: optimize  $\theta$  by minimizing the upper bound  $L$ .

## Analyzing $L$ , (1)

Goal: to derive a simplified objective based on  $L$ .

## Analyzing $L$ , (1)

Goal: to derive a simplified objective based on  $L$ .

Transition density of the forward process:

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

# Analyzing $L$ , (1)

Goal: to derive a simplified objective based on  $L$ .

Transition density of the forward process:

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k I_d).$$

## Lemma

For  $k > 0$ ,  $q(x^{(k)} | x^{(0)}) = \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k} x^{(0)}, (1 - \bar{\alpha}_k) I_d)$ , where

$$\alpha_k = 1 - \beta_k, \quad \text{and} \quad \bar{\alpha}_k = \prod_{i=1}^k \alpha_i. \tag{7}$$

# Analyzing $L$ , (1)

Goal: to derive a simplified objective based on  $L$ .

Transition density of the forward process:

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k I_d).$$

## Lemma

For  $k > 0$ ,  $q(x^{(k)} | x^{(0)}) = \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k} x^{(0)}, (1 - \bar{\alpha}_k) I_d)$ , where

$$\alpha_k = 1 - \beta_k, \quad \text{and} \quad \bar{\alpha}_k = \prod_{i=1}^k \alpha_i. \tag{7}$$

## Proof.

Prove by induction. □

## Analyzing $L$ , (2)

- ① Variational bound:  $L = \mathbb{E}_{\mathbb{Q}} \left( -\log p(x^{(N)}) - \sum_{k=1}^N \log \frac{p_{\theta}(x^{(k-1)} | x^{(k)})}{q(x^{(k)} | x^{(k-1)})} \right).$
- ② Bayes' theorem  $P(A|B) = P(B|A) \frac{P(A)}{P(B)}$

$$\implies q(x^{(k)} | x^{(k-1)}) = q(x^{(k-1)} | x^{(k)}, x^{(0)}) \frac{q(x^{(k)} | x^{(0)})}{q(x^{(k-1)} | x^{(0)})}.$$

## Analyzing $L$ , (2)

- ① Variational bound:  $L = \mathbb{E}_{\mathbb{Q}} \left( -\log p(x^{(N)}) - \sum_{k=1}^N \log \frac{p_{\theta}(x^{(k-1)} | x^{(k)})}{q(x^{(k)} | x^{(k-1)})} \right).$
- ② Bayes' theorem  $P(A|B) = P(B|A) \frac{P(A)}{P(B)}$

$$\implies q(x^{(k)} | x^{(k-1)}) = q(x^{(k-1)} | x^{(k)}, x^{(0)}) \frac{q(x^{(k)} | x^{(0)})}{q(x^{(k-1)} | x^{(0)})}.$$

- ③ Decomposition of  $L$ :

### Proposition

We have  $L = L_N + \sum_{k=2}^N L_{k-1} + L_0$ , where

$$L_0 = -\mathbb{E}_{\mathbb{Q}} \left( \log p_{\theta}(x^{(0)} | x^{(1)}) \right),$$

$$L_{k-1} = \mathbb{E}_{\mathbb{Q}} \left[ D_{KL} \left( q(x^{(k-1)} | x^{(k)}, x^{(0)}) \mid\mid p_{\theta}(x^{(k-1)} | x^{(k)}) \right) \right], \quad k = 2, \dots, N,$$

$$L_N = \mathbb{E}_{\mathbb{Q}} \left[ D_{KL} \left( q(x^{(N)} | x^{(0)}) \mid\mid p(x^{(N)}) \right) \right].$$

## Analyzing $L$ , (3)

Transition density

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

## Analyzing $L$ , (3)

Transition density

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

$$q(x^{(k)} | x^{(0)}) = \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k} x^{(0)}, (1 - \bar{\alpha}_k) \mathbf{I}_d),$$

## Analyzing $L$ , (3)

Transition density

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

$$q(x^{(k)} | x^{(0)}) = \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k} x^{(0)}, (1 - \bar{\alpha}_k) \mathbf{I}_d),$$

Bayes' theorem:

$$q(x^{(k-1)} | x^{(k)}, x^{(0)}) = q(x^{(k)} | x^{(k-1)}) \frac{q(x^{(k-1)} | x^{(0)})}{q(x^{(k)} | x^{(0)})}$$

## Analyzing $L$ , (3)

Transition density

$$q(x^{(k)} | x^{(k-1)}) = \mathcal{N}(x^{(k)}; \sqrt{1 - \beta_k} x^{(k-1)}, \beta_k \mathbf{I}_d).$$

$$q(x^{(k)} | x^{(0)}) = \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k} x^{(0)}, (1 - \bar{\alpha}_k) \mathbf{I}_d),$$

Bayes' theorem:

$$q(x^{(k-1)} | x^{(k)}, x^{(0)}) = q(x^{(k)} | x^{(k-1)}) \frac{q(x^{(k-1)} | x^{(0)})}{q(x^{(k)} | x^{(0)})}$$

Putting together,

### Lemma

We have  $q(x^{(k-1)} | x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k \mathbf{I}_d)$ , where

$$\begin{aligned}\tilde{\mu}_k(x^{(k)}, x^{(0)}) &= \frac{\sqrt{\bar{\alpha}_{k-1}} \beta_k}{1 - \bar{\alpha}_k} x^{(0)} + \frac{\sqrt{\alpha_k} (1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k} x^{(k)}, \\ \tilde{\beta}_k &= \frac{1 - \bar{\alpha}_{k-1}}{1 - \bar{\alpha}_k} \beta_k.\end{aligned}\tag{8}$$

## Analyzing $L$ , (4)

1  $q(x^{(k-1)} | x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k \mathbf{I}_d).$

## Analyzing $L$ , (4)

①  $q(x^{(k-1)} | x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k \mathbf{I}_d).$

② Choices of  $p_\theta(x^{(k-1)} | x^{(k)})$ :

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \sigma_k^2 \mathbf{I}_d\right),$$

where  $\sigma_k^2 = \beta_k$ , or  $\sigma_k^2 = \tilde{\beta}_k = \frac{1 - \bar{\alpha}_{k-1}}{1 - \bar{\alpha}_k} \beta_k$ .

## Analyzing $L$ , (4)

①  $q(x^{(k-1)} | x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k \mathbf{I}_d)$ .

② Choices of  $p_\theta(x^{(k-1)} | x^{(k)})$ :

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \sigma_k^2 \mathbf{I}_d\right),$$

where  $\sigma_k^2 = \beta_k$ , or  $\sigma_k^2 = \tilde{\beta}_k = \frac{1 - \bar{\alpha}_{k-1}}{1 - \bar{\alpha}_k} \beta_k$ .

### Lemma

For  $2 \leq k \leq N$ , we have

$$\begin{aligned} L_{k-1} &= \mathbb{E}_{\mathbb{Q}} \left[ D_{KL} \left( q(x^{(k-1)} | x^{(k)}, x^{(0)}) \mid\mid p_\theta(x^{(k-1)} | x^{(k)}) \right) \right] \\ &= \frac{1}{2\sigma_k^2} \mathbb{E}_{\mathbb{Q}} \left( |\mu_\theta(x^{(k)}, k) - \tilde{\mu}_k(x^{(k)}, x^{(0)})|^2 \right) + C. \end{aligned}$$

## Analyzing $L$ , (4)

①  $q(x^{(k-1)} | x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k \mathbf{I}_d).$

② Choices of  $p_\theta(x^{(k-1)} | x^{(k)})$ :

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \sigma_k^2 \mathbf{I}_d\right),$$

where  $\sigma_k^2 = \beta_k$ , or  $\sigma_k^2 = \tilde{\beta}_k = \frac{1 - \bar{\alpha}_{k-1}}{1 - \bar{\alpha}_k} \beta_k$ .

### Lemma

For  $2 \leq k \leq N$ , we have

$$\begin{aligned} L_{k-1} &= \mathbb{E}_{\mathbb{Q}} \left[ D_{KL} \left( q(x^{(k-1)} | x^{(k)}, x^{(0)}) \mid\mid p_\theta(x^{(k-1)} | x^{(k)}) \right) \right] \\ &= \frac{1}{2\sigma_k^2} \mathbb{E}_{\mathbb{Q}} \left( |\mu_\theta(x^{(k)}, k) - \tilde{\mu}_k(x^{(k)}, x^{(0)})|^2 \right) + C. \end{aligned}$$

Conclusion:  $\mu_\theta(x^{(k)}, k)$  should predict  $\tilde{\mu}_k$ .

## Analyzing $L$ , (5)

### ① Reparametrization:

$$\begin{aligned} q(x^{(k)}|x^{(0)}) &= \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k}x^{(0)}, (1 - \bar{\alpha}_k)\mathbf{I}_d) \\ \implies x^{(k)}(x^{(0)}, \epsilon) &= \sqrt{\bar{\alpha}_k}x^{(0)} + \sqrt{1 - \bar{\alpha}_k}\epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \mathbf{I}_d). \end{aligned}$$

②  $q(x^{(k-1)}|x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k\mathbf{I}_d)$ , where

$$\tilde{\mu}_k(x^{(k)}, x^{(0)}) = \frac{\sqrt{\bar{\alpha}_{k-1}}\beta_k}{1 - \bar{\alpha}_k}x^{(0)} + \frac{\sqrt{\alpha_k(1 - \bar{\alpha}_{k-1})}}{1 - \bar{\alpha}_k}x^{(k)}.$$

# Analyzing $L$ , (5)

## ① Reparametrization:

$$q(x^{(k)}|x^{(0)}) = \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k}x^{(0)}, (1 - \bar{\alpha}_k)\mathbf{I}_d)$$
$$\implies x^{(k)}(x^{(0)}, \epsilon) = \sqrt{\bar{\alpha}_k}x^{(0)} + \sqrt{1 - \bar{\alpha}_k}\epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \mathbf{I}_d).$$

②  $q(x^{(k-1)}|x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k\mathbf{I}_d)$ , where

$$\tilde{\mu}_k(x^{(k)}, x^{(0)}) = \frac{\sqrt{\bar{\alpha}_{k-1}}\beta_k}{1 - \bar{\alpha}_k}x^{(0)} + \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k}x^{(k)}.$$

$$L_{k-1} - C$$

$$= \mathbb{E}_{\mathbb{Q}} \left( \frac{1}{2\sigma_k^2} |\mu_\theta(x^{(k)}, k) - \tilde{\mu}_k(x^{(k)}, x^{(0)})|^2 \right)$$

$$= \mathbb{E}_{x^{(0)}, \epsilon} \left[ \frac{1}{2\sigma_k^2} \left| \mu_\theta(x^{(k)}(x^{(0)}, \epsilon), k) - \tilde{\mu}_k(x^{(k)}(x^{(0)}, \epsilon), \frac{1}{\sqrt{\bar{\alpha}_k}}(x^{(k)}(x^{(0)}, \epsilon) - \sqrt{1 - \bar{\alpha}_k}\epsilon)) \right|^2 \right]$$

$$= \mathbb{E}_{x^{(0)}, \epsilon} \left[ \frac{1}{2\sigma_k^2} \left| \mu_\theta(x^{(k)}(x^{(0)}, \epsilon), k) - \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)}(x^{(0)}, \epsilon) - \frac{\beta_k}{\sqrt{1 - \bar{\alpha}_k}}\epsilon \right) \right|^2 \right].$$

# Analyzing $L$ , (5)

## ① Reparametrization:

$$\begin{aligned} q(x^{(k)}|x^{(0)}) &= \mathcal{N}(x^{(k)}; \sqrt{\bar{\alpha}_k}x^{(0)}, (1 - \bar{\alpha}_k)\mathbf{I}_d) \\ \implies x^{(k)}(x^{(0)}, \epsilon) &= \sqrt{\bar{\alpha}_k}x^{(0)} + \sqrt{1 - \bar{\alpha}_k}\epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \mathbf{I}_d). \end{aligned}$$

②  $q(x^{(k-1)}|x^{(k)}, x^{(0)}) = \mathcal{N}(x^{(k-1)}; \tilde{\mu}_k(x^{(k)}, x^{(0)}), \tilde{\beta}_k\mathbf{I}_d)$ , where

$$\tilde{\mu}_k(x^{(k)}, x^{(0)}) = \frac{\sqrt{\bar{\alpha}_{k-1}}\beta_k}{1 - \bar{\alpha}_k}x^{(0)} + \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k}x^{(k)}.$$

$$L_{k-1} - C$$

$$= \mathbb{E}_{\mathbb{Q}} \left( \frac{1}{2\sigma_k^2} |\mu_\theta(x^{(k)}, k) - \tilde{\mu}_k(x^{(k)}, x^{(0)})|^2 \right)$$

$$= \mathbb{E}_{x^{(0)}, \epsilon} \left[ \frac{1}{2\sigma_k^2} \left| \mu_\theta(x^{(k)}(x^{(0)}, \epsilon), k) - \tilde{\mu}_k(x^{(k)}(x^{(0)}, \epsilon), \frac{1}{\sqrt{\bar{\alpha}_k}}(x^{(k)}(x^{(0)}, \epsilon) - \sqrt{1 - \bar{\alpha}_k}\epsilon)) \right|^2 \right]$$

$$= \mathbb{E}_{x^{(0)}, \epsilon} \left[ \frac{1}{2\sigma_k^2} \left| \mu_\theta(x^{(k)}(x^{(0)}, \epsilon), k) - \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)}(x^{(0)}, \epsilon) - \frac{\beta_k}{\sqrt{1 - \bar{\alpha}_k}}\epsilon \right) \right|^2 \right].$$

Hence, the parametrization:  $\mu_\theta(x^{(k)}, k) = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1 - \bar{\alpha}_k}}\epsilon_\theta(x^{(k)}, k) \right).$

## Analyzing $L$ , (6)

Transition density:

$$p_{\theta}(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_{\theta}(x^{(k)}, k), \sigma_k^2 \mathbf{I}_d\right),$$

with  $\mu_{\theta}(x^{(k)}, k) = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon_{\theta}(x^{(k)}, k) \right).$

## Analyzing $L$ , (6)

Transition density:

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \sigma_k^2 \mathbf{I}_d\right),$$

with  $\mu_\theta(x^{(k)}, k) = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon_\theta(x^{(k)}, k) \right).$

Sampling  $x^{(k-1)} \sim p_\theta(x^{(k-1)} | x^{(k)})$ :

$$x^{(k-1)} = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon_\theta(x^{(k)}, k) \right) + \sigma_k z, \text{ where } z \sim \mathcal{N}(0, \mathbf{I}_d). \quad (9)$$

## Analyzing $L$ , (6)

Transition density:

$$p_\theta(x^{(k-1)} | x^{(k)}) = \mathcal{N}\left(x^{(k-1)}; \mu_\theta(x^{(k)}, k), \sigma_k^2 \mathbf{I}_d\right),$$

with  $\mu_\theta(x^{(k)}, k) = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon_\theta(x^{(k)}, k) \right).$

Sampling  $x^{(k-1)} \sim p_\theta(x^{(k-1)} | x^{(k)})$ :

$$x^{(k-1)} = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon_\theta(x^{(k)}, k) \right) + \sigma_k z, \text{ where } z \sim \mathcal{N}(0, \mathbf{I}_d). \quad (9)$$

$$\begin{aligned} L_{k-1} - C &= \mathbb{E}_{x^{(0)}, \epsilon} \left[ \frac{1}{2\sigma_k^2} \left| \mu_\theta(x^{(k)}(x^{(0)}, \epsilon), k) - \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)}(x^{(0)}, \epsilon) - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon \right) \right|^2 \right] \\ &= \mathbb{E}_{x^{(0)}, \epsilon} \left[ \frac{\beta_k^2}{2\sigma_k^2 \alpha_k (1 - \bar{\alpha}_k)} \left| \epsilon_\theta(\sqrt{\bar{\alpha}_k} x^{(0)} + \sqrt{1 - \bar{\alpha}_k} \epsilon, k) - \epsilon \right|^2 \right]. \end{aligned}$$

## Loss and algorithm

Simplified loss:

$$L_{\text{simple}}(\theta) = \mathbb{E}_{k, x^{(0)}, \epsilon} \left[ \left| \epsilon_\theta(\sqrt{\bar{\alpha}_k} x^{(0)} + \sqrt{1 - \bar{\alpha}_k} \epsilon, k) - \epsilon \right|^2 \right]. \quad (10)$$

# Loss and algorithm

Simplified loss:

$$L_{\text{simple}}(\theta) = \mathbb{E}_{k, x^{(0)}, \epsilon} \left[ \left| \epsilon_\theta(\sqrt{\bar{\alpha}_k} x^{(0)} + \sqrt{1 - \bar{\alpha}_k} \epsilon, k) - \epsilon \right|^2 \right]. \quad (10)$$

---

## Algorithm Training

---

- 1: **repeat**
  - 2:      $x^{(0)} \sim q_0(x^{(0)})$
  - 3:      $k \sim \text{Uniform}(\{1, 2, \dots, N\})$
  - 4:      $\epsilon \sim \mathcal{N}(0, I_d)$ .
  - 5:     take gradient descent step on
  - 6:          $\nabla_\theta \left| \epsilon_\theta(\sqrt{\bar{\alpha}_k} x^{(0)} + \sqrt{1 - \bar{\alpha}_k} \epsilon, k) - \epsilon \right|^2$
  - 7: **until** converged
-

---

## Algorithm Sampling

---

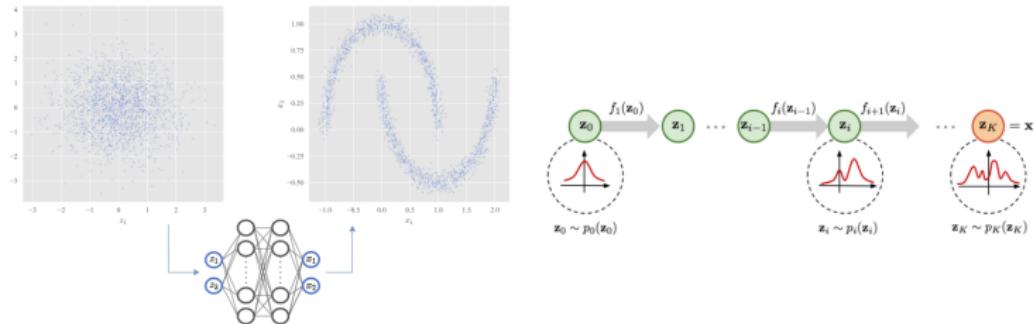
- 1:  $x^{(N)} \sim \mathcal{N}(0, \mathbf{I}_d)$
- 2: **for**  $k = N, \dots, 1$  **do**
- 3:      $z \sim \mathcal{N}(0, \mathbf{I}_d)$  if  $k > 1$ , **else**  $z = 0$
- 4:      $x^{(k-1)} = \frac{1}{\sqrt{\alpha_k}} \left( x^{(k)} - \frac{\beta_k}{\sqrt{1-\bar{\alpha}_k}} \epsilon_\theta(x^{(k)}, k) \right) + \sigma_k z$
- 5: **end for**
- 6: **return**  $x^{(0)}$

---

## Part 3: Normalizing flows

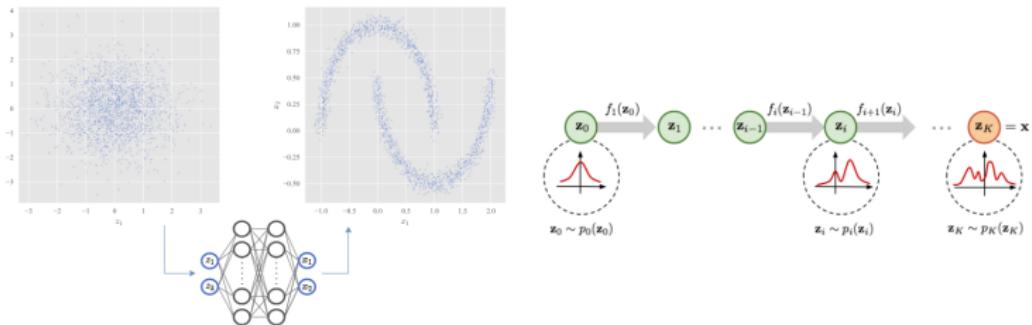
# Idea

Express data  $x \in \mathbb{R}^d$  as a transformation of (Gaussian)  $z \sim p_z(z)$ .



# Idea

Express data  $x \in \mathbb{R}^d$  as a transformation of (Gaussian)  $z \sim p_z(z)$ .



- ① “normalizing”:  $f^{-1}$  changes data distribution to normal distribution.
- ② “flow”:  $f$  is often composition of  $f_1, f_2, \dots, f_K$ .

## Change of variables

Let  $z \sim p_z(z)$  and  $x = f(z)$ , where  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is invertible and differentiable. Then,

$$p_x(x) = p_z(f^{-1}(x)) |\det J_{f^{-1}}(x)|. \quad (11)$$

## Change of variables

Let  $z \sim p_z(z)$  and  $x = f(z)$ , where  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is invertible and differentiable. Then,

$$p_x(x) = p_z(f^{-1}(x)) |\det J_{f^{-1}}(x)|. \quad (11)$$

- ① target  $p_x^*(x)$
- ② transformation  $x = f(z; \theta)$ .

Goal: to find  $\theta$  such that  $p_x(x; \theta)$  is close to  $p_x^*(x)$ .

## Loss function

$$p_x(x; \theta) = p_z(f^{-1}(x; \theta)) |\det J_{f^{-1}}(x; \theta)|.$$

# Loss function

$$p_x(x; \theta) = p_z(f^{-1}(x; \theta)) |\det J_{f^{-1}}(x; \theta)|.$$

(forward) KL divergence:

$$\begin{aligned} & D_{KL}\left(p_x^*(x) \mid p_x(x; \theta)\right) \\ &= \mathbb{E}_{x \sim p_x^*(x)}\left(\log \frac{p_x^*(x)}{p_x(x; \theta)}\right) \\ &= -\mathbb{E}_{x \sim p_x^*(x)}\left[\log p_z(f^{-1}(x; \theta)) + \log |\det J_{f^{-1}}(x; \theta)|\right] + C. \end{aligned}$$

# Loss function

$$p_x(x; \theta) = p_z(f^{-1}(x; \theta)) |\det J_{f^{-1}}(x; \theta)|.$$

(forward) KL divergence:

$$\begin{aligned} & D_{KL}\left(p_x^*(x) \mid p_x(x; \theta)\right) \\ &= \mathbb{E}_{x \sim p_x^*(x)}\left(\log \frac{p_x^*(x)}{p_x(x; \theta)}\right) \\ &= -\mathbb{E}_{x \sim p_x^*(x)}\left[\log p_z(f^{-1}(x; \theta)) + \log |\det J_{f^{-1}}(x; \theta)|\right] + C. \end{aligned}$$

Training objective in practice:

$$\text{Loss}(\theta) = -\frac{1}{N} \sum_{n=1}^N \left[ \log p_z(f^{-1}(x_n; \theta)) + \log |\det J_{f^{-1}}(x_n; \theta)| \right].$$

# Normalizing flow models

Requirements on  $f$ :

- ① input and output dimensions are the same;
- ②  $f$  is invertible;
- ③ computing  $f^{-1}$  and  $\det J_{f^{-1}}$  is efficient.

# Normalizing flow models

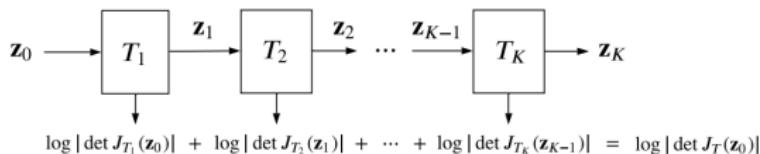
Requirements on  $f$ :

- ① input and output dimensions are the same;
- ②  $f$  is invertible;
- ③ computing  $f^{-1}$  and  $\det J_{f^{-1}}$  is efficient.

If  $f_1, f_2, \dots, f_K$  satisfy the above requirements, so does the composition  $f = f_K \circ \dots \circ f_1$ , and

$$\log |\det J_{f^{-1}}(x)| = \log \left| \prod_{k=1}^K \det J_{f_k^{-1}}(z_k) \right| = \sum_{k=1}^K \log |\det J_{f_k^{-1}}(z_k)|, \quad (12)$$

where  $z_K = x$  and  $z_{k-1} = f_k^{-1}(z_k)$ ,  $k = K, \dots, 1$ .



# Nonlinear Independent Components Estimation (NICE) model

$x = (x_1, x_2) = f(z)$  is defined as

$$\begin{aligned}x_1 &= z_1, \\x_2 &= z_2 + \mu_\theta(z_1),\end{aligned}$$

where  $z = (z_1, z_2)$ ,  $z_1 \in \mathbb{R}^{d'}$  and  $z_2 \in \mathbb{R}^{d-d'}$ , and  $\mu_\theta : \mathbb{R}^{d'} \rightarrow \mathbb{R}^{d-d'}$  is a neural network. The inverse of  $f$  is

$$\begin{aligned}z_1 &= x_1, \\z_2 &= x_2 - \mu_\theta(x_1).\end{aligned}$$

The Jacobian determinant of  $f$  is one.

# Real Non-Volume Preserving (RealNVP) model

$x = (x_1, x_2) = f(z)$  is defined as

$$x_1 = z_1,$$

$$x_2 = \exp(\sigma_\theta(z_1)) \odot z_2 + \mu_\theta(z_1),$$

where  $\odot$  denotes elementwise product, and  $\sigma_\theta, \mu_\theta : \mathbb{R}^{d'} \rightarrow \mathbb{R}^{d-d'}$ . The inverse of  $f$  is

$$z_1 = x_1,$$

$$z_2 = \exp(-\sigma_\theta(x_1)) \odot (x_2 - \mu_\theta(x_1)),$$

The Jacobian determinant is  $\det J_{f^{-1}}(x) = \exp(-\sum_{j=1}^{d-d'} (\sigma_\theta(x_1))_j)$ .

