

Denoising Diffusion Models: Fundamentals, Implementations and Applications

Mingfei Sun

University of Manchester

(Department of Computer Science)

July 1, 2024

Disclaimer

- ▶ Slides adapted from:

Denoising Diffusion Models: A Generative Learning Big Bang

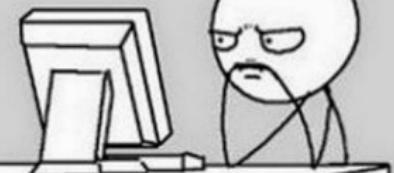
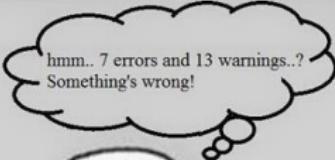
- ▶ Some images are taken from the blog:

What are Diffusion Models

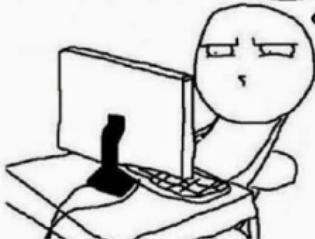
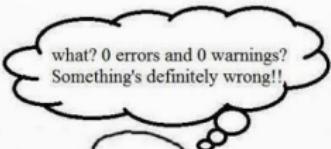
Heads-up: *n* warning(s) , 0 error(s)

Compiling program for the first time after coding

In second year



In final year



- ▶ This tutorial is math-intensive , containing many equations and derivations.
- ▶ This tutorial focuses specifically on de-noising diffusion models (DDPM) , with alterations in conditional generation.
- ▶ My knowledge on generative models is rather limited : may not answer all your questions, sorry!
- ▶ I tend to speak fast and skip many details : try to avoid it this time.

- ▶ This tutorial is math-intensive , containing many equations and derivations.
- ▶ This tutorial focuses specifically on de-noising diffusion models (DDPM) , with alterations in conditional generation.
- ▶ My knowledge on generative models is rather limited : may not answer all your questions, sorry!
- ▶ I tend to speak fast and skip many details : try to avoid it this time.

- ▶ This tutorial is math-intensive , containing many equations and derivations.
- ▶ This tutorial focuses specifically on de-noising diffusion models (DDPM) , with alterations in conditional generation.
- ▶ My knowledge on generative models is rather limited : may not answer all your questions, sorry!
- ▶ I tend to speak fast and skip many details : try to avoid it this time.

- ▶ This tutorial is math-intensive , containing many equations and derivations.
- ▶ This tutorial focuses specifically on de-noising diffusion models (DDPM) , with alterations in conditional generation.
- ▶ My knowledge on generative models is rather limited : may not answer all your questions, sorry!
- ▶ I tend to speak fast and skip many details : try to avoid it this time.

Outline

Reviews of Probability

Denoising Diffusion Probabilistic Models

Conditional Generation and Guidance

Applications beyond Image Generation

Outline

Reviews of Probability

Denoising Diffusion Probabilistic Models

Conditional Generation and Guidance

Applications beyond Image Generation

Random variables and notations

- ▶ A *random variable* (RV) is a *function* that assigns a number to the outcome of a random experiment.
- ▶ When referring to the probability $P(X = x)$, we usually simply write $P(x)$.
- ▶ Likewise, instead of writing $P(X = x, Y = y)$, we simply write $P(x, y)$.
- ▶ A continuous RV X takes values within one or more intervals of the real line.

Probability density function

- ▶ We use **probability density functions** (pdf), $p(x)$, to describe a continuous RV X .
- ▶ We can use a **joint probability density function**, $p(x, y)$ to fully characterise two continuous random variables X and Y .
- ▶ Typical pdf in diffusion models

$\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ with $\boldsymbol{\mu}$ as the mean and $\sigma^2 \mathbf{I}$ as the covariance matrix

- ▶ Normal distribution (aka Gaussian distribution) for scalar RV x :

$$\mathcal{N}(x; \boldsymbol{\mu}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \boldsymbol{\mu})^2}{2\sigma^2}\right)$$

- ▶ Multivariate normal distribution for vector RV \boldsymbol{x} :

$$\mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\boldsymbol{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\boldsymbol{x} - \boldsymbol{\mu})\right)$$

- ▶ For continuous RVs, expectation and variance are defined as

$$\mu = \mathbb{E}[X] = \int_{-\infty}^{\infty} xp(x)dx,$$

$$\sigma^2 = \text{var}[X] = \mathbb{E}[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 p(x)dx$$

Rules of probability (continuous RVs) and Markov property

- ▶ **Sum rule of probability.** In the case of continuous RVs, we replace the sums we had before with an integral

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy,$$

where $p(x)$ is known as the **marginal pdf**.

- ▶ **Product rule of probability.** The conditional pdf can be obtained as

$$p(x|y) = \frac{p(x, y)}{p(y)},$$

which can also be written as $p(x, y) = p(x|y)p(y)$.

- ▶ **Bayes' theorem** follows as

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)},$$

which can be verified through $p(x, y) = p(y|x)p(x)$ and $p(x, y) = p(x|y)p(y)$.

- ▶ **Markov property.** For all x_{t-1}, x_t, x_{t+1} :

$$p(x_{t+1}|x_t) = p(x_{t+1}|x_t, x_{t-1}, \dots, x_0),$$

i.e., x_t is a sufficient statistic for the history (x_0, x_1, \dots, x_t) .

Rules of probability (continuous RVs) and Markov property

- ▶ **Sum rule of probability.** In the case of continuous RVs, we replace the sums we had before with an integral

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy,$$

where $p(x)$ is known as the **marginal pdf**.

- ▶ **Product rule of probability.** The conditional pdf can be obtained as

$$p(x|y) = \frac{p(x, y)}{p(y)},$$

which can also be written as $p(x, y) = p(x|y)p(y)$.

- ▶ Bayes' theorem follows as

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)},$$

which can be verified through $p(x, y) = p(y|x)p(x)$ and $p(x, y) = p(x|y)p(y)$.

- ▶ **Markov property.** For all x_{t-1}, x_t, x_{t+1} :

$$p(x_{t+1}|x_t) = p(x_{t+1}|x_t, x_{t-1}, \dots, x_0),$$

i.e., x_t is a sufficient statistic for the history (x_0, x_1, \dots, x_t) .

Rules of probability (continuous RVs) and Markov property

- ▶ **Sum rule of probability.** In the case of continuous RVs, we replace the sums we had before with an integral

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy,$$

where $p(x)$ is known as the **marginal pdf**.

- ▶ **Product rule of probability.** The conditional pdf can be obtained as

$$p(x|y) = \frac{p(x, y)}{p(y)},$$

which can also be written as $p(x, y) = p(x|y)p(y)$.

- ▶ **Bayes' theorem** follows as

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)},$$

which can be verified through $p(x, y) = p(y|x)p(x)$ and $p(x, y) = p(x|y)p(y)$.

- ▶ **Markov property.** For all x_{t-1}, x_t, x_{t+1} :

$$p(x_{t+1}|x_t) = p(x_{t+1}|x_t, x_{t-1}, \dots, x_0),$$

i.e., x_t is a sufficient statistic for the history (x_0, x_1, \dots, x_t) .

Rules of probability (continuous RVs) and Markov property

- ▶ **Sum rule of probability.** In the case of continuous RVs, we replace the sums we had before with an integral

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy,$$

where $p(x)$ is known as the **marginal pdf**.

- ▶ **Product rule of probability.** The conditional pdf can be obtained as

$$p(x|y) = \frac{p(x, y)}{p(y)},$$

which can also be written as $p(x, y) = p(x|y)p(y)$.

- ▶ **Bayes' theorem** follows as

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)},$$

which can be verified through $p(x, y) = p(y|x)p(x)$ and $p(x, y) = p(x|y)p(y)$.

- ▶ **Markov property.** For all x_{t-1}, x_t, x_{t+1} :

$$p(x_{t+1}|x_t) = p(x_{t+1}|x_t, x_{t-1}, \dots, x_0),$$

i.e., x_t is a sufficient statistic for the history (x_0, x_1, \dots, x_t) .

Outline

Reviews of Probability

Denoising Diffusion Probabilistic Models

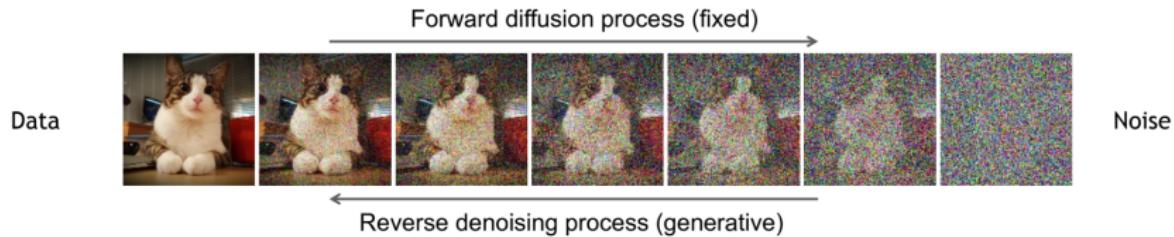
Conditional Generation and Guidance

Applications beyond Image Generation

Denoising Diffusion Models

Denoising Diffusion models consist of two processes:

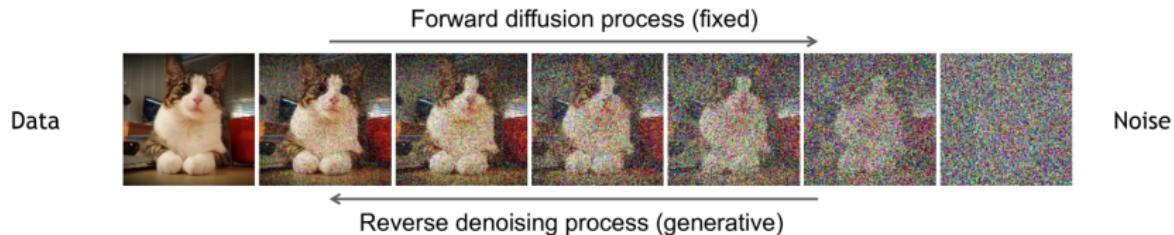
- ▶ Forward diffusion process that gradually adds noise to input
 - ▶ Reverse denoising process that learns to generate data by denoising



Denoising Diffusion Models

Denoising Diffusion models consist of two processes:

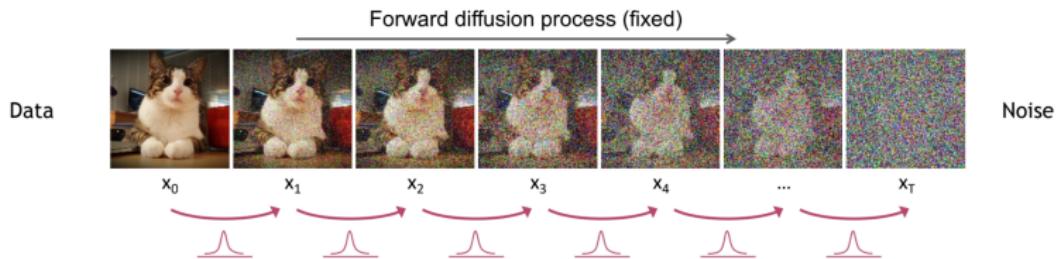
- ▶ Forward diffusion process that gradually adds noise to input
 - ▶ Reverse denoising process that learns to generate data by denoising



Forward Diffusion Models

The formal definition of the forward process in T steps:

- ▶ Add small amount of Gaussian noise to the sample in T steps, producing a sequence of noisy samples x_1, x_2, \dots, x_T .

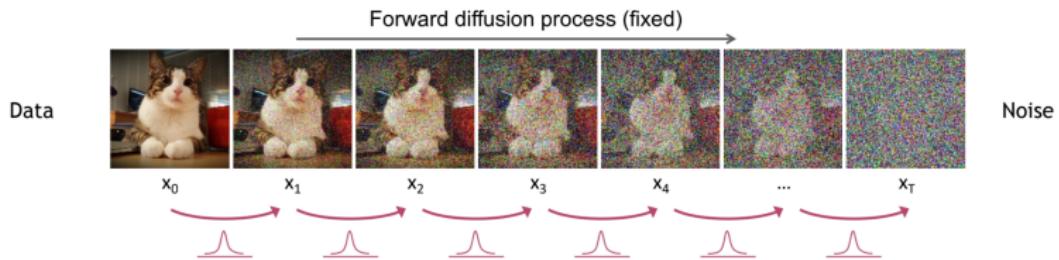


$$q(x_t | x_{t-1}) = \underbrace{\mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)}_{x_t \sim \mathcal{N}(\mu, \sigma^2) \text{ with } \mu = \sqrt{1 - \beta_t} x_{t-1} \text{ and } \sigma^2 I = \beta_t I}$$

Forward Diffusion Models

The formal definition of the forward process in T steps:

- ▶ Add small amount of Gaussian noise to the sample in T steps, producing a sequence of noisy samples x_1, x_2, \dots, x_T .



- ▶ The step sizes are controlled by a variance schedule $\{\beta_t \in (0, 1)\}_{t=1}^T$.

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \underbrace{\mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})}_{\mathbf{x}_t \sim \mathcal{N}(\mu, \sigma^2) \text{ with } \mu = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} \text{ and } \sigma^2 \mathbf{I} = \beta_t \mathbf{I}}$$

$$q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

= ...

$$= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon}$$

$$= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \epsilon && ; \text{ where } \epsilon \sim \mathcal{N}(0, 1) \\ &= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \epsilon && ; \text{ where } \epsilon \sim \mathcal{N}(0, 1) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon && ; \text{ where } \epsilon \sim \mathcal{N}(0, 1) \\ &= \dots \\ &= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \epsilon \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon && ; \text{ where } \epsilon \sim \mathcal{N}(0, 1) \end{aligned}$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(0, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$$

$$= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$$

= ...

$$= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon}$$

$$= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \dots \\ &= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon} \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})\end{aligned}$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \dots \\ &= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon} \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})\end{aligned}$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \dots \\ &= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon} \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})\end{aligned}$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \dots \\ &= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon} \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})\end{aligned}$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

Diffusion Kernel

A nice property of the forward process is that we can sample \mathbf{x}_t at any arbitrary timestep t in a closed form. Define $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.

According to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \mathcal{N}(\mathbf{0}, \beta_t \mathbf{I})$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t} \underbrace{\left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon} \right)}_{\mathbf{x}_{t-1}} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}) \\ &= \dots \\ &= \sqrt{\alpha_t \dots \alpha_0} \mathbf{x}_0 + \sqrt{1 - \alpha_t \dots \alpha_0} \boldsymbol{\epsilon} \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} && ; \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})\end{aligned}$$

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Or equivalently,

Diffusion Kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \mathcal{N}(\mathbf{0}, (1 - \bar{\alpha}_t) \mathbf{I})$

What happens to a distribution in the forward process

For sampling:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon},$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, β_t (a.k.a the noise scheduler) is designed such that $\bar{\alpha}_T \rightarrow 0$ and $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$.

So far, we discussed the diffusion kernel $q(\mathbf{x}_t | \mathbf{x}_0)$ but what about $q(\mathbf{x}_t)$?

$$\underbrace{q(\mathbf{x}_t)}_{\text{Diffused data distribution}} = \int \underbrace{q(\mathbf{x}_0, \mathbf{x}_t)}_{\text{Joint distribution}} d\mathbf{x}_0 = \int \underbrace{q(\mathbf{x}_0)}_{\text{Input data distribution}} \underbrace{q(\mathbf{x}_t | \mathbf{x}_0)}_{\text{Diffusion kernel}} d\mathbf{x}_0$$

We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)$ (i.e., ancestral sampling).

What happens to a distribution in the forward process

For sampling:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon},$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, β_t (a.k.a the noise scheduler) is designed such that $\bar{\alpha}_T \rightarrow 0$ and $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$.

So far, we discussed the diffusion kernel $q(\mathbf{x}_t | \mathbf{x}_0)$ but what about $q(\mathbf{x}_t)$?

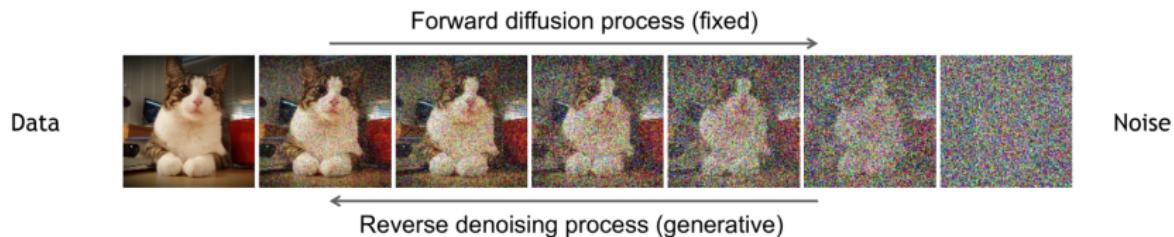
$$\underbrace{q(\mathbf{x}_t)}_{\text{Diffused data distribution}} = \int \underbrace{q(\mathbf{x}_0, \mathbf{x}_t)}_{\text{Joint distribution}} d\mathbf{x}_0 = \int \underbrace{q(\mathbf{x}_0)}_{\text{Input data distribution}} \underbrace{q(\mathbf{x}_t | \mathbf{x}_0)}_{\text{Diffusion kernel}} d\mathbf{x}_0$$

We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)$ (i.e., ancestral sampling).

Denoising Diffusion Models

Denoising Diffusion models consist of two processes:

- ▶ Forward diffusion process that gradually adds noise to input
- ▶ Reverse denoising process that learns to generate data by denoising



Generative Learning by Denoising

Heads-up: we use q for the forward process and p for the reverse process.

Recall, that the diffusion parameters are designed such that $p(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Generation:

- ▶ Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$
- ▶ Iteratively sample $\mathbf{x}_{t-1} \sim \underbrace{p(\mathbf{x}_{t-1} | \mathbf{x}_t)}_{\text{True denoising distribution}}$

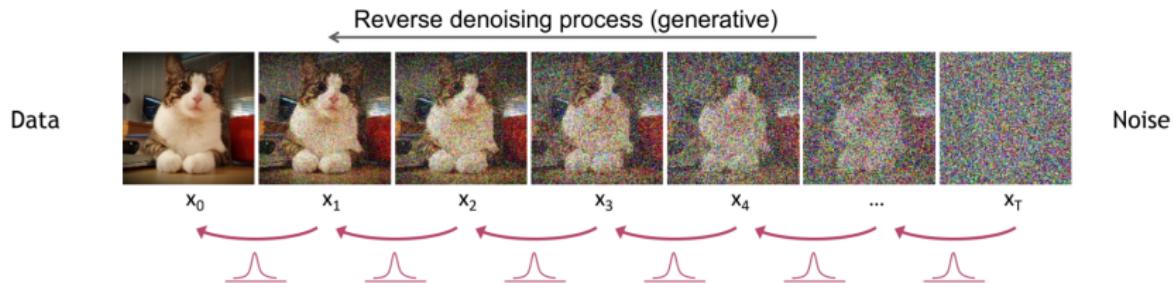
Can we parameterize $p(\mathbf{x}_{t-1} | \mathbf{x}_t)$? Yes, we can use a **Normal distribution** if β_t is small in each forward diffusion step.¹

¹For Gaussian diffusion, for continuous diffusion (limit of small step size β), the reversal of the diffusion process has the identical functional form as the forward process. Since $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ is a Gaussian distribution, and if β_t is small, then $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ will also be a Gaussian distribution.

Reverse Denoising Process

Formal definition of reverse processes in T steps:

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$
$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\mu_{\theta}(\mathbf{x}_t, t)}_{\text{Trainable network}} , \sigma_t^2 \mathbf{I})$$
$$\Rightarrow p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$



Learning Denoising Model: Variational upper bound

For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$\begin{aligned}\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] &\leq \mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] + \underbrace{\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]}_{\text{KL divergence}} \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right]}_{\triangleq \mathcal{L}}\end{aligned}$$

To derive a loss function for training $\mu_\theta(\mathbf{x}_t, t)$, with the following functions so far:

Forward process: $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$

Diffusion kernel: $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Reverse sampling from: $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Reverse transition: $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$

Learning Denoising Model: Variational upper bound

For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$\begin{aligned}\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] &\leq \mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] + \underbrace{\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]}_{\text{KL divergence}} \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right]}_{\triangleq \mathcal{L}}\end{aligned}$$

To derive a loss function for training $\mu_\theta(\mathbf{x}_t, t)$, with the following functions so far:

Forward process: $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$

Diffusion kernel: $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

Reverse sampling from: $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Reverse transition: $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$

Learning Denoising Model: Variational upper bound

$$\begin{aligned}\mathcal{L} &\triangleq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_{T-1})q(\mathbf{x}_{T-1}|\mathbf{x}_{T-2})...q(\mathbf{x}_1|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)p_{\theta}(\mathbf{x}_1|\mathbf{x}_2)...p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_T)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]\end{aligned}$$

Markov assumption of \mathbf{x}_t on \mathbf{x}_{t-1} :

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \textcolor{red}{\mathbf{x}_0})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

Apply Bayes' rule: $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)$

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right]$$

Learning Denoising Model: Variational upper bound

$$\begin{aligned}\mathcal{L} &\triangleq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_{T-1})q(\mathbf{x}_{T-1}|\mathbf{x}_{T-2})...q(\mathbf{x}_1|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)p_{\theta}(\mathbf{x}_1|\mathbf{x}_2)...p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_T)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]\end{aligned}$$

Markov assumption of \mathbf{x}_t on \mathbf{x}_{t-1} :

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \textcolor{red}{\mathbf{x}_0})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

Apply Bayes' rule: $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)$

$$\begin{aligned}&= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right]\end{aligned}$$

Learning Denoising Model: Variational upper bound

$$\begin{aligned}\mathcal{L} &\triangleq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_{T-1})q(\mathbf{x}_{T-1}|\mathbf{x}_{T-2})...q(\mathbf{x}_1|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)p_{\theta}(\mathbf{x}_1|\mathbf{x}_2)...p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_T)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]\end{aligned}$$

Markov assumption of \mathbf{x}_t on \mathbf{x}_{t-1} :

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \textcolor{red}{\mathbf{x}_0})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

Apply Bayes' rule: $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)$

$$\begin{aligned}&= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right]\end{aligned}$$

Learning Denoising Model: Variational upper bound

$$\begin{aligned}\mathcal{L} &\triangleq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_{T-1})q(\mathbf{x}_{T-1}|\mathbf{x}_{T-2})...q(\mathbf{x}_1|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)p_{\theta}(\mathbf{x}_1|\mathbf{x}_2)...p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_T)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]\end{aligned}$$

Markov assumption of \mathbf{x}_t on \mathbf{x}_{t-1} :

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \textcolor{red}{x_0})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

Apply Bayes' rule: $q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0) = q(x_{t-1}|x_t, x_0)q(x_t|x_0)$

$$\begin{aligned}&= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(x_{t-1}|x_t, x_0)}{p_{\theta}(x_{t-1}|x_t)} \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(x_{t-1}|x_t, x_0)}{p_{\theta}(x_{t-1}|x_t)} + \log \prod_{t=1}^T \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} \right]\end{aligned}$$

Learning Denoising Model: Variational upper bound

$$\begin{aligned}\mathcal{L} &\triangleq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_{T-1})q(\mathbf{x}_{T-1}|\mathbf{x}_{T-2})...q(\mathbf{x}_1|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)p_{\theta}(\mathbf{x}_1|\mathbf{x}_2)...p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_T)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]\end{aligned}$$

Markov assumption of \mathbf{x}_t on \mathbf{x}_{t-1} :

$$= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \textcolor{red}{\mathbf{x}_0})}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

Apply Bayes' rule: $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)$

$$\begin{aligned}&= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right]\end{aligned}$$

Learning Denoising Model: Variational upper bound

$$\begin{aligned}\mathcal{L} &\triangleq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{0:T})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log p(\mathbf{x}_T) + \log \prod_{t=1}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \prod_{t=1}^T \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right] \\&= \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\underbrace{\log p(\mathbf{x}_T)}_{\text{No learnable parameters}} + \boxed{\sum_{t=1}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}} + \underbrace{\log q(\mathbf{x}_T|\mathbf{x}_0)}_{\text{No learnable parameters}} \right]\end{aligned}$$

Learning Denoising Model: Variational upper bound

For training, we can form variational upper bound

$$\nabla_{\theta} \mathcal{L} = \nabla_{\theta} \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is tractable and follows Gaussian distribution (shown later):

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} = \frac{\overbrace{q(\mathbf{x}_t|\mathbf{x}_{t-1})}^{\text{Forward process}} \overbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}^{\text{Diffusion kernel}}}{\underbrace{q(\mathbf{x}_t|\mathbf{x}_0)}_{\text{Diffusion kernel}}}$$

Recall that

- ▶ Forward process:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) \sim \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

- ▶ Diffusion kernel:

$$\begin{aligned} q(\mathbf{x}_{t-1}|\mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0, (1 - \bar{\alpha}_{t-1}) \mathbf{I}), \\ q(\mathbf{x}_t|\mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \end{aligned}$$

- ▶ Reverse transition

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

Learning Denoising Model: Variational upper bound

For training, we can form variational upper bound

$$\nabla_{\theta} \mathcal{L} = \nabla_{\theta} \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is tractable and follows Gaussian distribution (shown later):

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} = \frac{\overbrace{q(\mathbf{x}_t|\mathbf{x}_{t-1})}^{\text{Forward process}} \overbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}^{\text{Diffusion kernel}}}{\underbrace{q(\mathbf{x}_t|\mathbf{x}_0)}_{\text{Diffusion kernel}}}$$

Recall that

- ▶ Forward process:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) \sim \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

- ▶ Diffusion kernel:

$$\begin{aligned} q(\mathbf{x}_{t-1}|\mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0, (1 - \bar{\alpha}_{t-1}) \mathbf{I}), \\ q(\mathbf{x}_t|\mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \end{aligned}$$

- ▶ Reverse transition

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

Learning Denoising Model

$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ follows Gaussian distribution $\mathcal{N}\left(\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}\right)$:

$$\begin{aligned} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \\ &\propto q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) q(\mathbf{x}_{t-1} | \mathbf{x}_0) \\ &\propto \exp\left(-\frac{(\mathbf{x}_t - \sqrt{1-\beta_t} \mathbf{x}_{t-1})^2}{\beta_t} - \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0)^2}{1-\bar{\alpha}_{t-1}}\right) \\ &\propto \exp\left(-\frac{\mathbf{x}_t^2 - 2\sqrt{1-\beta_t} \mathbf{x}_t \mathbf{x}_{t-1} + \alpha_t \mathbf{x}_{t-1}^2}{\beta_t} - \frac{\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_{t-1} \mathbf{x}_0 + \bar{\alpha}_{t-1} \mathbf{x}_0^2}{1-\bar{\alpha}_{t-1}}\right) \\ &\propto \exp\left[-\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right) \mathbf{x}_{t-1}^2 + \left(\frac{2\sqrt{1-\beta_t}}{\beta_t} \mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}} \mathbf{x}_0\right) \mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0)\right] \end{aligned}$$

So the mean is given by

$$\begin{aligned} \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) &= \left(\frac{\sqrt{1-\beta_t}}{\beta_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}} \mathbf{x}_0\right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right) \\ &= \frac{(1-\bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1-\bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \mathbf{x}_0 \end{aligned}$$

Recall that $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t} \boldsymbol{\epsilon}$. Substitute $\mathbf{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t} \boldsymbol{\epsilon})$

$$\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon} \right)$$

Learning Denoising Model: Variational upper bound

For training, we can form variational upper bound

$$\nabla_{\theta} \mathcal{L} = \nabla_{\theta} \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ follows Gaussian distribution $\mathcal{N}\left(\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}\right)$, where
 $\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon} \right)$.

Since \mathbf{x}_t is available as input to the model, we may choose the parameterization

$$\boxed{\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right)}$$

Using a few simple arithmetic operations and $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t} \boldsymbol{\epsilon}$, we can write down the variational objective for two Gaussian distributions as:

$$\min_{\theta} \mathbb{E}_{\substack{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim \mathcal{U}\{1, T\}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \text{data distribution} \quad \text{Diffusion steps} \quad \text{Gaussian noises}}} \left[\lambda_t \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2 \right]$$

Ho et al.² observe that simply setting λ_t to 1 for all t works best in practice.

²Denoising diffusion probabilistic models, Ho et al, 2020

Learning Denoising Model: Variational upper bound

For training, we can form variational upper bound

$$\nabla_{\theta} \mathcal{L} = \nabla_{\theta} \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ follows Gaussian distribution $\mathcal{N}\left(\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}\right)$, where
 $\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon} \right)$.

Since \mathbf{x}_t is available as input to the model, we may choose the parameterization

$$\boxed{\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right)}$$

Using a few simple arithmetic operations and $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t} \boldsymbol{\epsilon}$, we can write down the variational objective for two Gaussian distributions as:

$$\min_{\theta} \mathbb{E}_{\substack{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim \mathcal{U}\{1, T\}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \text{data distribution} \quad \text{Diffusion steps} \quad \text{Gaussian noises}}} \left[\lambda_t \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2 \right]$$

Ho et al.² observe that simply setting λ_t to 1 for all t works best in practice.

²Denoising diffusion probabilistic models, Ho et al, 2020

Generating New Samples

Denoising steps:

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

and

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right)$$

Note that $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, a new sample \mathbf{x}_0 is generated by recursively applying

$$\begin{aligned} \mathbf{x}_{t-1} &= \mu_{\theta}(\mathbf{x}_t, t) + \sigma_t \mathbf{z} & ; \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \forall t = T, \dots, 1 \\ &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z} & ; \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \forall t = T, \dots, 1 \end{aligned}$$

- ▶ Ho et al.³ suggest to set σ_t to $\sigma_t^2 = \beta_t$ or $\sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$, yielding similar experimental results.
- ▶ Nichol and Dhariwal⁴ propose to learn an interpolation between β_t and $\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$

³Denoising diffusion probabilistic models, Ho et al, 2020

⁴Improved denoising diffusion probabilistic models, Alexander Quinn Nichol and Prafulla Dhariwal, 2021

Combined: Training and Sampling

Practical algorithms for Denoising Diffusion Probabilistic Models (DDPM):

Algorithm Training

```
1: for  $i = 1$  to  $N$  do
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(1, \dots, T)$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on  $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$ 
6: end for
```

Algorithm Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( \mathbf{x}_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5:   return  $\mathbf{x}_0$ 
6: end for
```

Combined: Training and Sampling

Practical algorithms for Denoising Diffusion Probabilistic Models (DDPM):

Algorithm Training

```
1: for  $i = 1$  to  $N$  do
2:    $x_0 \sim q(x_0)$ 
3:    $t \sim \text{Uniform}(1, \dots, T)$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on  $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$ 
6: end for
```

Algorithm Sampling

```
1:  $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $x_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right) + \sigma_t z$ 
5:   return  $x_0$ 
6: end for
```

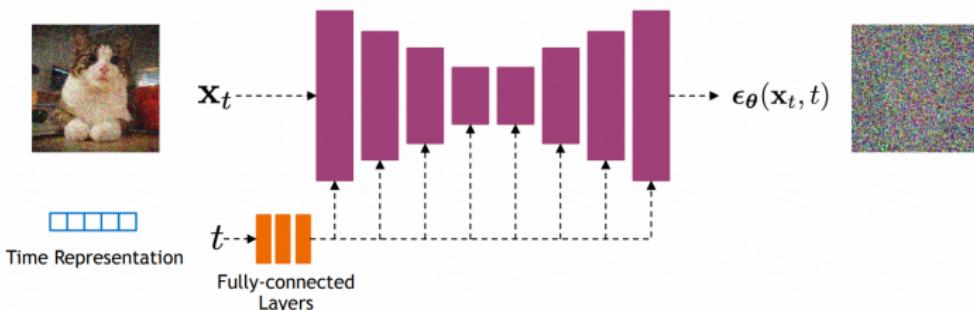
Implementation considerations

Algorithm Training

```
1: for  $i = 1$  to  $N$  do
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(1, \dots, T)$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on  $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2$ 
6: end for
```

Network architectures for ϵ_{θ} :

- ▶ Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent $\epsilon_{\theta}(\mathbf{x}_t, t)$
- ▶ Time representation: sinusoidal positional embeddings or random Fourier features



Implementation considerations

Algorithm Training

```
1: for  $i = 1$  to  $N$  do
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(1, \dots, T)$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on  $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta} \left( \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2$ 
6: end for
```

Algorithm Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5:   return  $\mathbf{x}_0$ 
6: end for
```

Noise and variance scheduler:

- ▶ Set $T = 1000$ and the forward process variances to constants increasing linearly from $\beta_1 = 10^{-4}$ to $\beta_T = 0.02$.
- ▶ $\alpha_t \triangleq 1 - \beta_t$ and $\bar{\alpha}_t \triangleq \prod_{i=1}^t \alpha_i$.
- ▶ Diagonal covariance matrix and $\sigma_t^2 = \beta_t$.

Outline

Reviews of Probability

Denoising Diffusion Probabilistic Models

Conditional Generation and Guidance

Applications beyond Image Generation

Impressive Conditional Diffusion Models

Text-to-image generation:

- ▶ Text inputs as conditions
- ▶ Training conditional models



panda mad scientist mixing sparkling chemicals, artstation



A cute corgi lives in a house made out of sushi.

Conditioning and Guidance Techniques

A conditional diffusion model is a modification of an unconditional diffusion model, by conditioning on a RV \mathbf{y} :

$$\text{Unconditional: } p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$\text{Conditional: } p_{\theta}(\mathbf{x}_{0:T} | \mathbf{y}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$$

► Classifier guidance

- To apply Bayes' rule to factorize the conditions as an extra classifier model to guide the sampling process

► Explicit conditions

• A condition can be explicitly specified as a function of the input image, such as a mask or a set of coordinates.

► Classifier-free guidance

• A condition can be implicitly specified as a function of the input image, such as a mask or a set of coordinates.

Conditioning and Guidance Techniques

A conditional diffusion model is a modification of an unconditional diffusion model, by conditioning on a RV \mathbf{y} :

$$\text{Unconditional: } p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$\text{Conditional: } p_{\theta}(\mathbf{x}_{0:T} | \mathbf{y}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$$

- ▶ **Classifier guidance**
 - ▶ To apply Bayes' rule to factorize the conditions as an extra classifier model to guide the sampling process
- ▶ **Explicit conditions**
 - ▶ To feed conditions as an explicit input to the noise predictor in both training and sampling
- ▶ **Classifier-free guidance**
 - ▶ To apply Bayes' rule to factorize the conditions and then jointly train a conditional and unconditional diffusion model through dropout

Conditioning and Guidance Techniques

A conditional diffusion model is a modification of an unconditional diffusion model, by conditioning on a RV \mathbf{y} :

$$\text{Unconditional: } p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$\text{Conditional: } p_{\theta}(\mathbf{x}_{0:T} | \mathbf{y}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$$

► Classifier guidance

- To apply Bayes' rule to factorize the conditions as an extra classifier model to guide the sampling process

► Explicit conditions

- To feed conditions as an explicit input to the noise predictor in both training and sampling

► Classifier-free guidance

- To apply Bayes' rule to factorize the conditions and then jointly train a conditional and unconditional diffusion model through dropout

Conditioning and Guidance Techniques

A conditional diffusion model is a modification of an unconditional diffusion model, by conditioning on a RV \mathbf{y} :

$$\text{Unconditional: } p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$\text{Conditional: } p_{\theta}(\mathbf{x}_{0:T} | \mathbf{y}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$$

- ▶ **Classifier guidance**
 - ▶ To apply Bayes' rule to factorize the conditions as an extra classifier model to guide the sampling process
- ▶ **Explicit conditions**
 - ▶ To feed conditions as an explicit input to the noise predictor in both training and sampling
- ▶ **Classifier-free guidance**
 - ▶ To apply Bayes' rule to factorize the conditions and then jointly train a conditional and unconditional diffusion model through dropout

Conditioning and Guidance Techniques

A conditional diffusion model is a modification of an unconditional diffusion model, by conditioning on a RV \mathbf{y} :

$$\text{Unconditional: } p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$\text{Conditional: } p_{\theta}(\mathbf{x}_{0:T} | \mathbf{y}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$$

- ▶ **Classifier guidance**
 - ▶ To apply Bayes' rule to factorize the conditions as an extra classifier model to guide the sampling process
- ▶ **Explicit conditions**
 - ▶ To feed conditions as an explicit input to the noise predictor in both training and sampling
- ▶ **Classifier-free guidance**
 - ▶ To apply Bayes' rule to factorize the conditions and then jointly train a conditional and unconditional diffusion model through dropout

Conditioning and Guidance Techniques

A conditional diffusion model is a modification of an unconditional diffusion model, by conditioning on a RV \mathbf{y} :

$$\text{Unconditional: } p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$\text{Conditional: } p_{\theta}(\mathbf{x}_{0:T} | \mathbf{y}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$$

- ▶ **Classifier guidance**
 - ▶ To apply Bayes' rule to factorize the conditions as an extra classifier model to guide the sampling process
- ▶ **Explicit conditions**
 - ▶ To feed conditions as an explicit input to the noise predictor in both training and sampling
- ▶ **Classifier-free guidance**
 - ▶ To apply Bayes' rule to factorize the conditions and then jointly train a conditional and unconditional diffusion model through dropout

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

To condition the reverse de-noising process $p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})$ on a label \mathbf{y} , We can apply Bayes' rule to sample each transition according to

$$p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) \propto p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) \underbrace{p(\mathbf{y}|\mathbf{x}_t)}_{\text{Classifier}}$$

Recall that the diffusion model predicts the previous timestep \mathbf{x}_t from timestep \mathbf{x}_{t+1} using a Gaussian distribution:

$$\begin{aligned} p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) &= \mathcal{N}(\mu, \Sigma) \\ \Rightarrow \log p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) &= -\frac{1}{2}(\mathbf{x}_t - \mu)^\top \Sigma (\mathbf{x}_t - \mu) \end{aligned}$$

Now we parameterize $p(\mathbf{y}|\mathbf{x}_t)$ with ϕ and approximate $\log p_\phi(\mathbf{y}|\mathbf{x}_t)$ using a Taylor expansion around $\mathbf{x}_t = \mu$: $\log p_\phi(\mathbf{y}|\mathbf{x}_t) \approx (\mathbf{x}_t - \mu)^\top \underbrace{\nabla_{\mathbf{x}_t} p_\phi(\mathbf{y}|\mathbf{x}_t)}_{\text{Define as } g} + C_1$. Hence,

$$\begin{aligned} \log p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) &\approx -\frac{1}{2}(\mathbf{x}_t - \mu)^\top \Sigma (\mathbf{x}_t - \mu) + (\mathbf{x}_t - \mu)^\top g + C_2 \\ &\approx -\frac{1}{2}(\mathbf{x}_t - \mu - \Sigma g)^\top \Sigma (\mathbf{x}_t - \mu - \Sigma g) + g^\top \Sigma g + C_2 \\ &\approx \underbrace{-\frac{1}{2}(\mathbf{x}_t - \mu - \Sigma g)^\top \Sigma (\mathbf{x}_t - \mu - \Sigma g) + C_3}_{=\log \mathcal{N}(\mu + \Sigma g, \Sigma)} \end{aligned}$$

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

To condition the reverse de-noising process $p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})$ on a label \mathbf{y} , We can apply Bayes' rule to sample each transition according to

$$p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) \propto p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) \underbrace{p(\mathbf{y}|\mathbf{x}_t)}_{\text{Classifier}}$$

Recall that the diffusion model predicts the previous timestep \mathbf{x}_t from timestep \mathbf{x}_{t+1} using a Gaussian distribution:

$$\begin{aligned} p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) &= \mathcal{N}(\boldsymbol{\mu}, \Sigma) \\ \Rightarrow \log p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) &= -\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu})^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu}) \end{aligned}$$

Now we parameterize $p(\mathbf{y}|\mathbf{x}_t)$ with ϕ and approximate $\log p_\phi(\mathbf{y}|\mathbf{x}_t)$ using a Taylor expansion around $\mathbf{x}_t = \boldsymbol{\mu}$: $\log p_\phi(\mathbf{y}|\mathbf{x}_t) \approx (\mathbf{x}_t - \boldsymbol{\mu})^\top \underbrace{\nabla_{\mathbf{x}_t} p_\phi(\mathbf{y}|\mathbf{x}_t)}_{\text{Define as } g} + C_1$. Hence,

$$\begin{aligned} \log p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) &\approx -\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu})^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu}) + (\mathbf{x}_t - \boldsymbol{\mu})^\top g + C_2 \\ &\approx -\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu} - \Sigma g)^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu} - \Sigma g) + g^\top \Sigma g + C_2 \\ &\approx \underbrace{-\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu} - \Sigma g)^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu} - \Sigma g) + C_3}_{=\log \mathcal{N}(\boldsymbol{\mu} + \Sigma g, \Sigma)} \end{aligned}$$

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

To condition the reverse de-noising process $p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})$ on a label \mathbf{y} , We can apply Bayes' rule to sample each transition according to

$$p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) \propto p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) \underbrace{p(\mathbf{y}|\mathbf{x}_t)}_{\text{Classifier}}$$

Recall that the diffusion model predicts the previous timestep \mathbf{x}_t from timestep \mathbf{x}_{t+1} using a Gaussian distribution:

$$\begin{aligned} p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) &= \mathcal{N}(\boldsymbol{\mu}, \Sigma) \\ \Rightarrow \log p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}) &= -\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu})^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu}) \end{aligned}$$

Now we parameterize $p(\mathbf{y}|\mathbf{x}_t)$ with ϕ and approximate $\log p_\phi(\mathbf{y}|\mathbf{x}_t)$ using a Taylor expansion around $\mathbf{x}_t = \boldsymbol{\mu}$: $\log p_\phi(\mathbf{y}|\mathbf{x}_t) \approx (\mathbf{x}_t - \boldsymbol{\mu})^\top \underbrace{\nabla_{\mathbf{x}_t} p_\phi(\mathbf{y}|\mathbf{x}_t)}_{\text{Define as } \mathbf{g}} + C_1$. Hence,

$$\begin{aligned} \log p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) &\approx -\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu})^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu}) + (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{g} + C_2 \\ &\approx -\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu} - \Sigma \mathbf{g})^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu} - \Sigma \mathbf{g}) + \mathbf{g}^\top \Sigma \mathbf{g} + C_2 \\ &\approx \underbrace{-\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu} - \Sigma \mathbf{g})^\top \Sigma (\mathbf{x}_t - \boldsymbol{\mu} - \Sigma \mathbf{g}) + C_3}_{=\log \mathcal{N}(\boldsymbol{\mu} + \Sigma \mathbf{g}, \Sigma)} \end{aligned}$$

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

The conditional transition operator can be approximated by a Gaussian similar to the unconditional transition operator, but with its mean shifted by $\Sigma \mathbf{g}$.

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) \propto \mathcal{N}(\boldsymbol{\mu} + \Sigma \mathbf{g}, \Sigma)$$

where $\mathbf{g} \triangleq \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t)$ is the gradient of a classifier.

How to obtain such gradients?

- ▶ Train a classifier p_ϕ with the diffused samples and their labels tuples, i.e., $(\mathbf{x}_t, \mathbf{y})$ where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$.
- ▶ Compute the log-likelihood gradients of this classifier w.r.t input \mathbf{x}_t , NOT ϕ .
- ▶ Add a scale s to classifier gradients to trade off sample fidelity against diversity.
- ▶ Some implementations consider the step as an extra input: $p_\phi(\mathbf{y} | \mathbf{x}_t, t)$.

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

The conditional transition operator can be approximated by a Gaussian similar to the unconditional transition operator, but with its mean shifted by Σg .

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) \propto \mathcal{N}(\boldsymbol{\mu} + \Sigma \mathbf{g}, \Sigma)$$

where $\mathbf{g} \triangleq \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t)$ is the gradient of a classifier.

How to obtain such gradients?

- ▶ Train a classifier p_ϕ with the diffused samples and their labels tuples, i.e., $(\mathbf{x}_t, \mathbf{y})$ where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$.
- ▶ Compute the log-likelihood gradients of this classifier w.r.t input \mathbf{x}_t , NOT ϕ .
- ▶ Add a scale s to classifier gradients to trade off sample fidelity against diversity.
- ▶ Some implementations consider the step as an extra input: $p_\phi(y | \mathbf{x}_t, \mathbf{t})$.

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

The conditional transition operator can be approximated by a Gaussian similar to the unconditional transition operator, but with its mean shifted by $\Sigma \mathbf{g}$.

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) \propto \mathcal{N}(\boldsymbol{\mu} + \Sigma \mathbf{g}, \Sigma)$$

where $\mathbf{g} \triangleq \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t)$ is the gradient of a classifier.

How to obtain such gradients?

- ▶ Train a classifier p_ϕ with the diffused samples and their labels tuples, i.e., $(\mathbf{x}_t, \mathbf{y})$ where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$.
- ▶ Compute the log-likelihood gradients of this classifier w.r.t input \mathbf{x}_t , NOT ϕ .
- ▶ Add a scale s to classifier gradients to trade off sample fidelity against diversity.
- ▶ Some implementations consider the step as an extra input: $p_\phi(y | \mathbf{x}_t, \mathbf{t})$.

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

The conditional transition operator can be approximated by a Gaussian similar to the unconditional transition operator, but with its mean shifted by Σg .

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) \propto \mathcal{N}(\boldsymbol{\mu} + \Sigma \mathbf{g}, \Sigma)$$

where $\mathbf{g} \triangleq \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t)$ is the gradient of a classifier.

How to obtain such gradients?

- ▶ Train a classifier p_ϕ with the diffused samples and their labels tuples, i.e., $(\mathbf{x}_t, \mathbf{y})$ where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$.
- ▶ Compute the log-likelihood gradients of this classifier w.r.t input \mathbf{x}_t , NOT ϕ .
- ▶ Add a scale s to classifier gradients to trade off sample fidelity against diversity.
- ▶ Some implementations consider the step as an extra input: $p_\phi(\mathbf{y} | \mathbf{x}_t, \textcolor{red}{t})$.

Classifier Guidance: Bayes' Rule in Action

Exploiting class labels to guide sampling process

Algorithm Classifier training (new)

- 1: **for** $i = 1$ to N **do**
- 2: $(\mathbf{x}_0, \mathbf{y}) \sim q(\mathbf{x}_0, \mathbf{y})$
- 3: $t \sim \text{Uniform}(1, \dots, T)$
- 4: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$
- 5: Take gradient descent step on $-\nabla_{\phi} \log p_{\phi}(\mathbf{y}|\mathbf{x}_t)$
- 6: **end for**

Algorithm Guided sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(\mathbf{x}_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + s \cdot \sigma_t \nabla_{\mathbf{x}_t} \log p_{\phi}(\mathbf{y}|\mathbf{x}_t) + \sigma_t \mathbf{z}$
- 5: **return** \mathbf{x}_0
- 6: **end for**

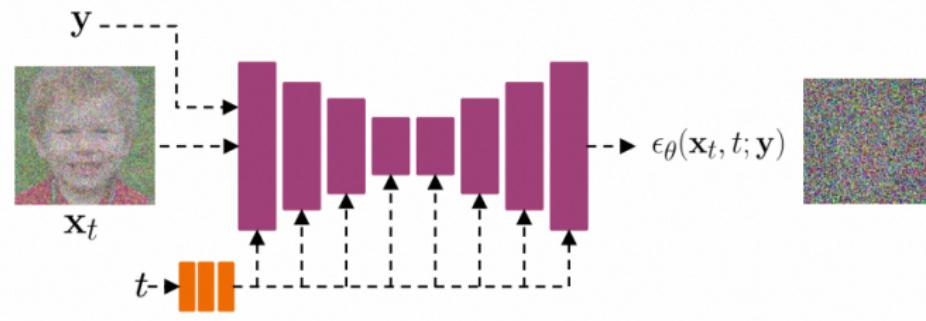
Explicit Conditional Training

Conditional sampling can be considered as training $p(\mathbf{x}|\mathbf{y})$ where \mathbf{y} is the input conditioning (e.g., text) and \mathbf{x} is generated output (e.g., image).

Train the score model for \mathbf{x} conditioned on \mathbf{y} using:

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}(\mathbf{x}, \mathbf{y}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}\{0, T\}} \|\epsilon_\theta(\mathbf{x}_t, t; \mathbf{y}) - \epsilon\|^2$$

The conditional score is simply a U-Net with \mathbf{x}_t and \mathbf{y} together in the input.



Explicit Conditional DDPM

Algorithm Training

```
1: for  $i = 1$  to  $N$  do
2:    $(\mathbf{x}_0, \mathbf{y}) \sim q(\mathbf{x}_0, \mathbf{y})$ 
3:    $t \sim \text{Uniform}(1, \dots, T)$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on  $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \mathbf{y}, t) \right\|^2$ 
6: end for
```

Algorithm Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: Set  $\mathbf{y}$ 
3: for  $t = T, \dots, 1$  do
4:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, \mathbf{y}, t) \right) + \sigma_t \mathbf{z}$ 
6: return  $\mathbf{x}_0$ 
7: end for
```

Classifier-free Guidance

Recall that classifier guidance requires training a classifier.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log \frac{p(\mathbf{y} | \mathbf{x}_t) p(\mathbf{x}_t)}{p(\mathbf{y})} = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)}_{\text{Classifier gradient}} + \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Score model}}$$

Re-arrange the equation:

$$\underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)}_{\text{Classifier gradient}} = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})}_{\text{Conditional diffusion model}} - \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Unconditional diffusion model}}$$

- ▶ Instead of training an additional classifier, get an “implicit classifier” by jointly training a conditional and unconditional diffusion model.
- ▶ In practice, the conditional and unconditional models are trained together by randomly dropping the condition of the diffusion model at certain chance.

Introduce a guidance scale ($1 + w$):

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) &\approx (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - w \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\end{aligned}$$

Classifier-free Guidance

Recall that classifier guidance requires training a classifier.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log \frac{p(\mathbf{y} | \mathbf{x}_t) p(\mathbf{x}_t)}{p(\mathbf{y})} = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)}_{\text{Classifier gradient}} + \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Score model}}$$

Re-arrange the equation:

$$\underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)}_{\text{Classifier gradient}} = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})}_{\text{Conditional diffusion model}} - \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Unconditional diffusion model}}$$

- ▶ Instead of training an additional classifier, get an “implicit classifier” by jointly training a conditional and unconditional diffusion model.
- ▶ In practice, the conditional and unconditional models are trained together by randomly dropping the condition of the diffusion model at certain chance.

Introduce a guidance scale ($1 + w$):

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) &\approx (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - w \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\end{aligned}$$

Classifier-free Guidance

Recall that classifier guidance requires training a classifier.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log \frac{p(\mathbf{y} | \mathbf{x}_t) p(\mathbf{x}_t)}{p(\mathbf{y})} = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)}_{\text{Classifier gradient}} + \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Score model}}$$

Re-arrange the equation:

$$\underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)}_{\text{Classifier gradient}} = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})}_{\text{Conditional diffusion model}} - \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Unconditional diffusion model}}$$

- ▶ Instead of training an additional classifier, get an “implicit classifier” by jointly training a conditional and unconditional diffusion model.
- ▶ In practice, the conditional and unconditional models are trained together by randomly dropping the condition of the diffusion model at certain chance.

Introduce a guidance scale ($1 + w$):

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) &\approx (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - w \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\end{aligned}$$

Classifier-free Diffusion Guidance

Algorithm Training

```
1: for  $i = 1$  to  $N$  do
2:    $(\mathbf{x}_0, \mathbf{y}) \sim q(\mathbf{x}_0, \mathbf{y})$ 
3:    $t \sim \text{Uniform}(1, \dots, T)$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   mask = Bernoulli( $p$ )
6:   Take gradient descent step on  $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \mathbf{y} * \text{mask}, t) \right\|^2$ 
7: end for
```

Algorithm Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: Set  $\mathbf{y}$ , guidance scale  $w$ 
3: for  $t = T, \dots, 1$  do
4:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:    $\epsilon^{(1)} = \epsilon_{\theta}(\mathbf{x}_t, \mathbf{y} * 1, t)$ 
6:    $\epsilon^{(0)} = \epsilon_{\theta}(\mathbf{x}_t, \mathbf{y} * 0, t)$ 
7:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left[ \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} ((1 + w)\epsilon^{(1)} - w\epsilon^{(0)}) \right] + \sigma_t \mathbf{z}$ 
8:   return  $\mathbf{x}_0$ 
9: end for
```

Give me the code, please!

Talk is cheap, show me the code!!!

Give me the code, please!

Educational implementations of DDPM, Guided DDPM, Conditional DDPM, and Classifier Free Diffusion, on MNIST dataset, i.e., generating digits



Outline

Reviews of Probability

Denoising Diffusion Probabilistic Models

Conditional Generation and Guidance

Applications beyond Image Generation

Applications beyond Image Generation

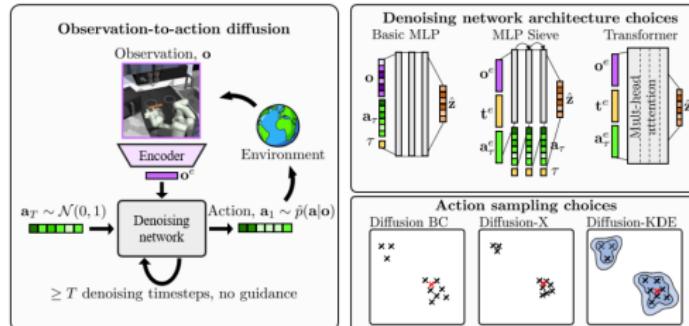
Diffusion models are good at modeling complex distributions

- ▶ Policies in reinforcement learning and imitation learning are conditional distributions:
 - ▶ Diffusion models as implicit parameterization of policies
 - ▶ Diffusion models to generate decision-making rules
 - ▶ Example: Imitating Human Behaviour with Diffusion Models, ICLR 2023
- ▶ Inference under constraints can be interpreted as constraints-guided sampling
 - ▶ Inference without constraints as generation modeling
 - ▶ Guide the sampling process with constraints-induced classifiers
 - ▶ Example: Effective Generation of Feasible Solutions for Integer Programming via Guided Diffusion, SIGKDD 2024

Diffusion Models for Imitation Learning

Explicit conditional diffusion models as policy representations⁵

- ▶ Policies as observation-to-action generative models



- ▶ Denoising Diffusion Probabilistic Models

$$\mathcal{L}_{\text{DDPM}} \triangleq \mathbb{E}_{o, a, t, z} [\| \epsilon_{\theta}(o, a_t, t) - z \|]$$

with $a_{\tau} = \sqrt{\bar{\alpha}_t}a + \sqrt{1 - \bar{\alpha}_t}z$ for variance scheduler $\bar{\alpha}_t$, and denosing

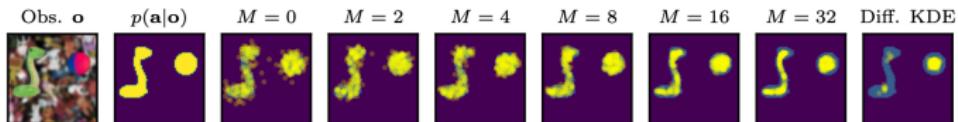
$$a_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(a_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \underbrace{\epsilon(o, a_t, t)}_{\text{Explicit conditioning}} \right) + \sigma_t z$$

⁵Imitating Human Behaviour with Diffusion Models, ICLR 2023

Diffusion Models for Imitation Learning

Explicit conditional diffusion models as policy representations⁶

- ▶ “A bad action could be selected during a roll-out”
- ▶ Reliable sampling schemes:
 - ▶ Diffusion-X: extra denoising iterations for a sample
 - ▶ Diffusion-KDE: kernel-density estimator (KDE) fitting



- ▶ Imitating human behavior [Counter-Strike: Global Offensive (CSGO)]



⁶Imitating Human Behaviour with Diffusion Models, ICLR 2023

Diffusion Models for Inference under Constraints

Guided sampling as inference under constraints⁷

Integer Programming

$$\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \quad \mathbf{x} \in \mathbb{Z}^n$$

where $\mathbf{c} \in \mathbb{R}^n$ denotes the objective coefficient, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the coefficient matrix of constraints and $\mathbf{b} \in \mathbb{R}^m$ represents the right-hand-side vector.

- ▶ Solving large-scale Integer Programming needs an initial good guess of feasible solutions
- ▶ Introduce **constraint guidance** by designing each transition probability as

$$p_{\theta, \phi}(\mathbf{z}_x^{(t)} | \mathbf{z}_x^{(t+1)}, \mathbf{z}_i, \mathbf{A}, \mathbf{b}) = Z p_\theta(\mathbf{z}_x^{(t)} | \mathbf{z}_x^{(t+1)}, \mathbf{z}_i) e^{-sc_\phi(\mathbf{z}_x^{(t)}, \mathbf{z}_i, \mathbf{A}, \mathbf{b})}, \quad (1)$$

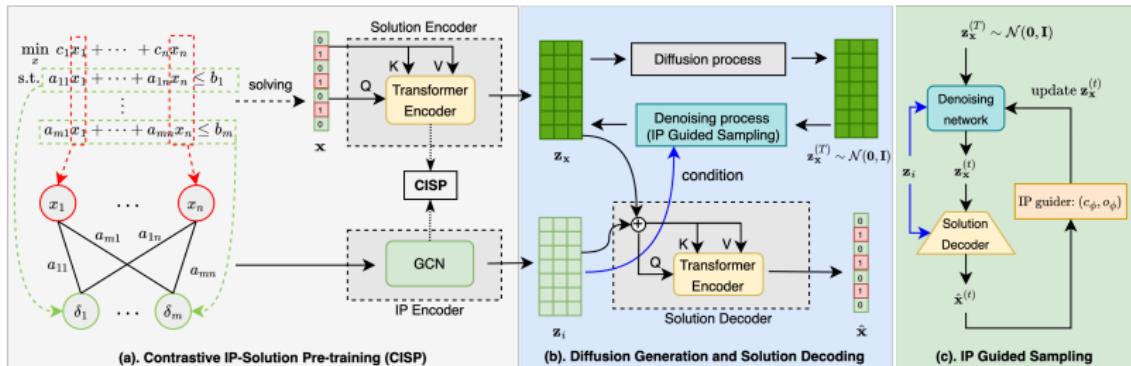
where $c_\phi(\mathbf{z}_x^{(t)}, \mathbf{z}_i, \mathbf{A}, \mathbf{b}) = \sum_{k=1}^m \max(\mathbf{a}_k^T \mathbf{d}_\phi(\mathbf{z}_x^{(t)}, \mathbf{z}_i) - b_k, 0)$ measures the violation of constraints.

- ▶ Consider the Taylor expansion for c_ϕ at $\mathbf{z}_x^{(t)} = \mu$:
 $c_\phi(\mathbf{z}_x^{(t)}, \mathbf{z}_i, \mathbf{A}, \mathbf{b}) = (\mathbf{z}_x^{(t)} - \mu) \nabla_{\mathbf{z}_x^{(t)}} c_\phi(\mathbf{z}_x^{(t)}, \mathbf{z}_i, \mathbf{A}, \mathbf{b})|_{\mathbf{z}_x^{(t)}=\mu} + C_1$. Hence, train a "classifier" to predict the constraint violation.

⁷Effective Generation of Feasible Solutions for Integer Programming via Guided Diffusion, SIGKDD 2024

Diffusion Models for Integer Programming

Guided sampling as inference under constraints⁸



1. Trains IP Encoder and Solution Encoder to obtain embeddings.
2. Jointly train diffusion models and the solution decoder to capture solution distributions.
3. Guide diffusion sampling with objectives and constraints.

⁸Effective Generation of Feasible Solutions for Integer Programming via Guided Diffusion, SIGKDD 2024

Conclusions

Thank You



mingfei.sun@manchester.ac.uk

<https://mingfeisun.github.io/>