

Introduction to Diffusion Models

MIAO Zhou

The Hong Kong Polytechnic University

January 6, 2025



- ▶ **Introduction**
- ▶ Preliminary Knowledge
- ▶ DDPM DDIM
 - ▶ Score based diffusion model
 - ▶ SDE diffusion
 - ▶ Reference

Brief Introduction to Diffusion Models



Diffusion models are a class of generative models that progressively transform simple noise into complex data distributions through a series of learned denoising steps.

- These models are inspired by physical diffusion processes, where data is iteratively corrupted by noise and then reconstructed.
- They have shown state-of-the-art performance in image generation, speech synthesis, and other generative tasks.

In this presentation, we will cover the **basic concepts** of typical models and **illustration of code**.



- ▶ Introduction
- ▶ Preliminary Knowledge
- ▶ DDPM DDIM
 - ▶ Score based diffusion model
 - ▶ SDE diffusion
 - ▶ Reference

Definition of ELBO



Generative Model Objective: The primary goal of generative models is to approximate the data distribution $p(\mathbf{x})$, allowing us to model and generate data. The optimization objective is to maximize the log-likelihood of the observed data:

$$\log p(\mathbf{x})$$

Challenge: Directly optimizing $\log p(\mathbf{x})$ is often intractable due to the complexity of marginalizing over latent variables \mathbf{z} :

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}.$$

Solution: Evidence Lower Bound (ELBO): To overcome this challenge, we introduce a variational approximation $q(\mathbf{z} | \mathbf{x})$ to the true posterior and derive a lower bound:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\frac{\log p(\mathbf{x}, \mathbf{z})}{\log q_{\phi}(\mathbf{z} | \mathbf{x})} \right].$$

The ELBO can be rewritten as:

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) \parallel p(\mathbf{z})),$$

where:

- $\log p(\mathbf{x} | \mathbf{z})$: Reconstruction term, measuring how well the latent variable explains the data.
- $\text{KL}(q\|p)$: Regularization term, ensuring $q(\mathbf{z} | \mathbf{x})$ is close to the prior $p(\mathbf{z})$.

Comparing VAE and Diffusion Models via ELBO



Variational Autoencoders (VAE):

- VAE optimizes the Evidence Lower Bound (ELBO) of the log-likelihood:

$$\text{ELBO}(\mathbf{x}) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log \overbrace{p_{\theta}(\mathbf{x}|\mathbf{z})}^{\text{a Gaussian}}]}_{\text{how good your decoder is}} - \underbrace{\mathbb{D}_{\text{KL}}(\overbrace{q_{\phi}(\mathbf{z}|\mathbf{x})}^{\text{a Gaussian}} \parallel \overbrace{p(\mathbf{z})}^{\text{a Gaussian}})}_{\text{how good your encoder is}}.$$

Diffusion Models:

- Diffusion models also optimize a variational lower bound but frame the process as a series of gradual denoising steps:

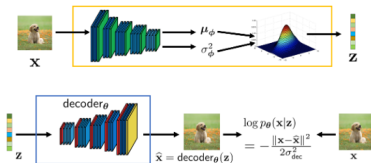
$$\begin{aligned} \text{ELBO}_{\phi, \theta}(\mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{x}_1|\mathbf{x}_0)}[\log \underbrace{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}_{\text{how good the initial block is}}] - \mathbb{E}_{q_{\phi}(\mathbf{x}_{T-1}|\mathbf{x}_0)}[\underbrace{\mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))}_{\text{how good the final block is}}] \\ &\quad - \sum_{t=1}^{T-1} \mathbb{E}_{q_{\phi}(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)}[\underbrace{\mathbb{D}_{\text{KL}}(q_{\phi}(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))}_{\text{how good the transition blocks are}}]. \end{aligned}$$

Architectures of VAE and Diffusion Model



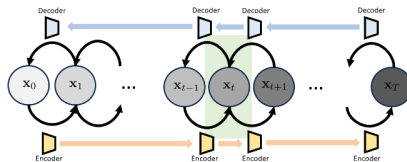
Variational Autoencoders (VAE):

- VAE Architecture:



Diffusion Models:

- Diffusion Model Architecture:





► Introduction

► Preliminary Knowledge

► **DDPM DDIM**

► Score based diffusion model

► SDE diffusion

► Reference



Introduction to Parameters in DDPM Model



Key Parameters in DDPM:

- T : Total number of timesteps in the diffusion process.
 - Controls the granularity of the diffusion process.
 - A larger T allows for more gradual noise addition and removal but increases computational cost.
- β_t : Noise schedule parameter for each timestep t .
 - Represents the variance of the added noise at timestep t .
 - Typically increases linearly or quadratically over time to ensure smooth noise corruption.
 - Values are bounded: $0 < \beta_t < 1$.
- α_t : Derived from β_t as $\alpha_t = 1 - \beta_t$.
 - Represents the retained signal at timestep t .
 - $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ accumulates the retained signal across timesteps, helping to express intermediate distributions.

Formulas of forward and backward process



Conditional distribution $q_\phi(\mathbf{x}_t \mid \mathbf{x}_0)$. The conditional distribution $q_\phi(\mathbf{x}_t \mid \mathbf{x}_0)$ is given by

$$q_\phi(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t \mid \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

where $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$.

The distribution $q_\phi(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$ takes the form of

$$q_\phi(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1} \mid \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0), \boldsymbol{\Sigma}_q(t))$$

where

$$\begin{aligned} \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) &= \frac{(1 - \bar{\alpha}_{t-1}) \sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t) \sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_0 \\ \boldsymbol{\Sigma}_q(t) &= \frac{(1 - \alpha_t)(1 - \sqrt{\bar{\alpha}_{t-1}})}{1 - \bar{\alpha}_t} \mathbf{I} \stackrel{\text{def}}{=} \sigma_q^2(t) \mathbf{I}. \end{aligned}$$

Training and Inference of DDPM



The loss function for a denoising diffusion probabilistic model:

$$\theta^* = \operatorname{argmin}_{\theta} \sum_{t=1}^T \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t) - \mathbf{x}_0\|^2 \right]$$

- **Training algorithm**

- Pick a random time stamp $\ell \sim \text{Uniform}[1, T]$.
- Draw a sample $\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t | \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$, i.e.,

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}).$$

- Take the gradient $\nabla_{\theta} \|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t) - \mathbf{x}_0\|^2$ to update $\hat{\mathbf{x}}_{\theta}$

Training and Inference of DDPM



- Inference on a Denoising Diffusion Probabilistic Model. (Version: Predict image)
- You give us a white noise vector $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$.
- Repeat the following for $t = T, T-1, \dots, 1$.
- We calculate $\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t)$ using our trained denoiser.
- Update according to

$$\mathbf{x}_{t-1} = \frac{(1 - \bar{\alpha}_{t-1}) \sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t) \sqrt{\alpha_{t-1}}}{1 - \bar{\alpha}_t} \hat{\mathbf{x}}_{\theta}(\mathbf{x}_t) + \sigma_q(t) \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$$



Let us consider a family \mathcal{Q} of inference distributions, indexed by a real vector $\sigma \in \mathbb{R}_{\geq 0}^T$:

$$q_{\sigma}(x_{1:T} \mid x_0) := q_{\sigma}(x_T \mid x_0) \prod_{t=2}^T q_{\sigma}(x_{t-1} \mid x_t, x_0)$$

where $q_{\sigma}(x_T \mid x_0) = \mathcal{N}(\sqrt{\alpha_T}x_0, (1 - \alpha_T)\mathbf{I})$ and for all $t > 1$,

$$q_{\sigma}(x_{t-1} \mid x_t, x_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}x_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{x_t - \sqrt{\alpha_t}x_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right).$$

The mean function is chosen to order to ensure that $q_{\sigma}(x_t \mid x_0) = \mathcal{N}(\sqrt{\alpha_t}x_0, (1 - \alpha_t)\mathbf{I})$ for all t (see Lemma 1 of Appendix B), so that it defines a joint inference distribution that matches the "marginals" as desired. The forward process³ can be derived from Bayes' rule:

$$q_{\sigma}(x_t \mid x_{t-1}, x_0) = \frac{q_{\sigma}(x_{t-1} \mid x_t, x_0) q_{\sigma}(x_t \mid x_0)}{q_{\sigma}(x_{t-1} \mid x_0)}$$

Inference of DDIM



Training function of DDIM is the same as DDPM, sampling function one can generate a sample \mathbf{x}_{t-1} from a sample \mathbf{x}_t via:

$$x_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left(\frac{x_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}^{(t)}(x_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } \mathbf{x}_0"} + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2}}_{\text{"direction pointing to } \mathbf{x}_t"} \cdot \epsilon_{\theta}^{(t)}(x_t) + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

where $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is standard Gaussian noise independent of \mathbf{x}_t , and we define $\alpha_0 := 1$. Different choices of σ values results in different generative processes, all while using the same model ϵ_{θ} , so re-training the model is unnecessary. When $\sigma_t = \sqrt{(1 - \alpha_{t-1}) / (1 - \alpha_t)} \sqrt{1 - \alpha_t / \alpha_{t-1}}$ for all t , the forward process becomes Markovian, and the generative process becomes a DDPM.



- ▶ Introduction
- ▶ Preliminary Knowledge
- ▶ DDPM DDIM
 - ▶ **Score based diffusion model**
 - ▶ SDE diffusion
 - ▶ Reference

Introduction to Langevin Dynamics



The Langevin dynamics for sampling from a known distribution $p(\mathbf{x})$ is an iterative procedure for $t = 1, \dots, T$:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \tau \nabla_{\mathbf{x}} \log p(\mathbf{x}_t) + \sqrt{2\tau} \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$$

where τ is the step size which users can control, and \mathbf{x}_0 is white noise.

Without the noise term $\sqrt{2\tau} \mathbf{z}$ at the end, the Langevin dynamics is literally gradient descent. The descent direction $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ is carefully chosen that \mathbf{x}_t will converge to the distribution $p(\mathbf{x})$.



- **Loss Function of Exact Score Matching (ESM):**

$$\min_{\theta} \mathcal{J}_{\text{ESM}}(\theta) := \mathbb{E}_{p(t, X)} \left[\|s_{\theta}(t, X) - \nabla \log p(t, X)\|^2 \right]$$

which expands to:

$$\mathcal{J}_{\text{ESM}}(\theta) = \|s_{\theta}(t, X)\|^2 - 2\nabla \cdot s_{\theta}(t, X)$$

- **Loss Function of Sliced Score Matching (SSM):**

$$\min_{\theta} \tilde{\mathcal{J}}_{\text{SSM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{v \sim \mathcal{N}(0, I)} \mathbb{E}_{p(t, X)} \left[\lambda(t) \left(\|s_{\theta}(t, X)\|^2 + 2v^{\top} \nabla(v^{\top} s_{\theta}(t, X)) \right) \right]$$

Loss Function of Score Matching



- **Loss Function of Denoising Score Matching(DSM):**

$$\min J_{\text{DSM}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathbb{E}_{q(\mathbf{x}, \mathbf{x}')} \left[\frac{1}{2} \left\| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}) - \nabla_{\mathbf{x}} q(\mathbf{x} | \mathbf{x}') \right\|^2 \right]$$

which expands to:

$$\mathbb{E}_{q(\mathbf{x}')} \left[\frac{1}{2} \left\| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}' + \sigma \mathbf{z}) + \frac{\mathbf{z}}{\sigma} \right\|^2 \right]$$



- ▶ Introduction
- ▶ Preliminary Knowledge
- ▶ DDPM DDIM
 - ▶ Score based diffusion model
 - ▶ SDE diffusion
 - ▶ Reference

Training and Inference of SDE Diffusion



- Forward Diffusion.

$$d\mathbf{x} = \underbrace{\mathbf{f}(\mathbf{x}, t)}_{\text{drift}} dt + \underbrace{g(t)}_{\text{diffusion}} d\mathbf{w}.$$

- Reverse Diffusion.

$$d\mathbf{x} = \underbrace{\mathbf{f}(\mathbf{x}, t)}_{\text{drift}} - g(t)^2 \underbrace{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}_{\text{score function}} dt + \underbrace{g(t)d\bar{\mathbf{w}}}_{\text{reverse-time diffusion}}$$

- Forward Diffusion function of DDPM.

$$d\mathbf{x} = \underbrace{-\frac{\beta(t)}{2}\mathbf{x}}_{=\mathbf{f}(\mathbf{x}, t)} dt + \underbrace{\sqrt{\beta(t)} d\mathbf{w}}_{=g(t)}$$

- Reverse Diffusion function of DDPM.

$$d\mathbf{x} = -\beta(t) \left[\frac{\mathbf{x}}{2} + \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}$$



- ▶ Introduction
- ▶ Preliminary Knowledge
- ▶ DDPM DDIM
 - ▶ Score based diffusion model
 - ▶ SDE diffusion
- ▶ Reference



- **Articles:**

- Tutorial on Diffusion Models for Imaging and Vision [arxiv](#).
- SCORE-BASED DIFFUSION MODELS VIA STOCHASTIC DIFFERENTIAL EQUATIONS – A TECHNICAL TUTORIAL [arxiv](#).

- **Github Repositories:**

- Useful project of DDPM and DDIM based on pytorch [Github](#).
- Useful project of score matching diffusion based on pytorch [Github](#).
- All illustration code and PPT pdf [Github](#).

- **Youtubers:**

- Dhiraj Madan [Youtube](#).

Good Luck!



- Thank you for your attention!
- We hope this introduction has been helpful for understanding the basics.
- Wishing you success in your journey with diffusion models!



THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學

Opening Minds • Shaping the Future • 啟迪思維 • 成就未來