

CS5242 : Neural Networks and Deep Learning

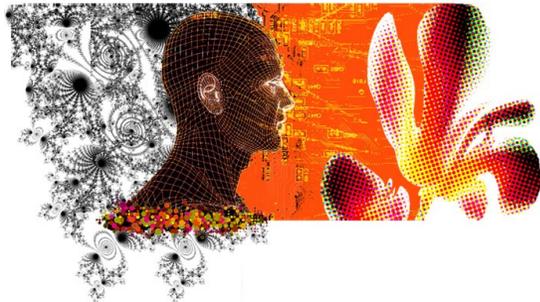
Lecture 8 : Diffusion Models

Semester 2 2024/25

Xavier Bresson

<https://x.com/xbresson>

Department of Computer Science
National University of Singapore (NUS)



Outline

- Introduction
- Vanilla diffusion models (DDPM)
- Lower bound on data distribution
- Noising process
- Denoising process
- Learning to denoise
- Conclusion

Outline

- **Introduction**
- Vanilla diffusion models (DDPM)
- Lower bound on data distribution
- Noising process
- Denoising process
- Learning to denoise
- Conclusion

Introduction

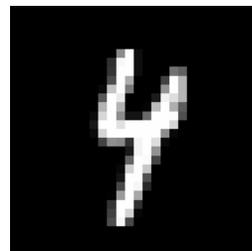
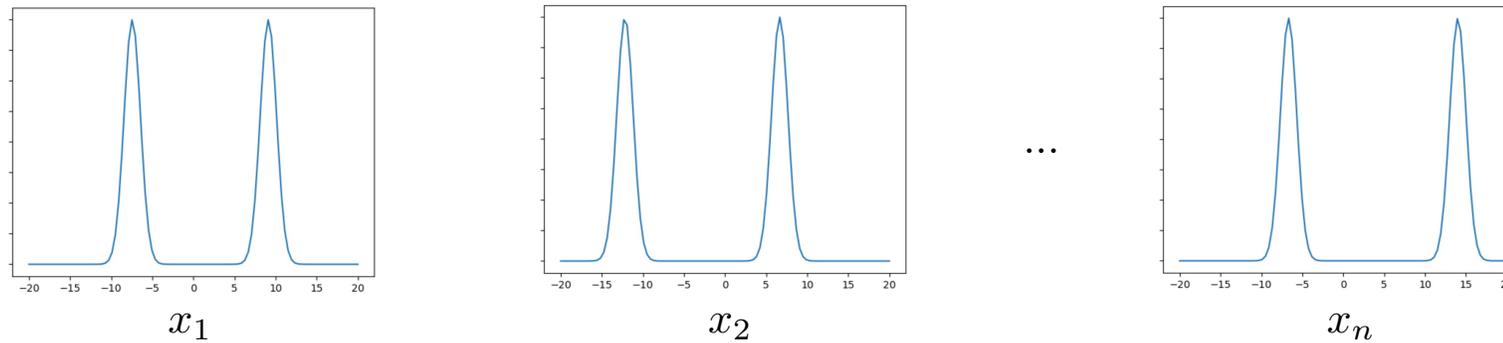
- Generative models can be broadly categorized based on the type of data they generate -- discrete or continuous.
- For discrete data (e.g. a sequence of words, where each word belongs to a finite dictionary of tokens), the most common models are auto-regressive models. These include Recurrent Neural Networks (RNNs) and Transformer-based Language Models (LMs).
- For continuous data (e.g. an image represented as a grid of pixels with continuous values in a RGB color space $[0,255]^3$), the primary generative models include:
- Variational Autoencoders (VAEs): Stable and robust during training, but often produce less accurate results.
- Generative Adversarial Networks (GANs): Capable of generating highly realistic outputs, but difficult to train.
- Diffusion Models (DMs): Combine robustness in training with high-quality outputs. DMs have become the state-of-the-art for image generation in computer vision.

Lecture approach

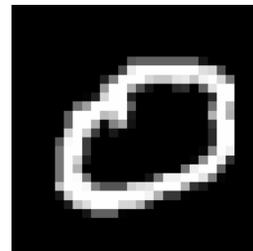
- In this lecture, we focus on Diffusion Models.
- Our main goal is to derive the governing equations of diffusion models from fundamental statistical principles -- using only basic probabilistic tools such as Bayes' theorem, the expectation of function of random variables, etc.
- The lecture is self-contained.
- While the full derivation is lengthy, each step is sound and requires no additional computations beyond what is presented.

Datasets

- We will use the following datasets :
 - An artificial dataset of mixtures of two identical Gaussian distributions with random shifts. The data dimensionality is $d_x = N$ (all sequences have N continuous variables).
 - The MNIST image dataset with dimensionality $d_x = N_x \times N_y \times 3$ (N_x pixel width, N_y pixel height, and each pixel has 3 color values).

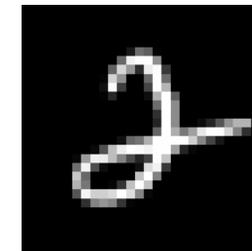


x_1



x_2

...



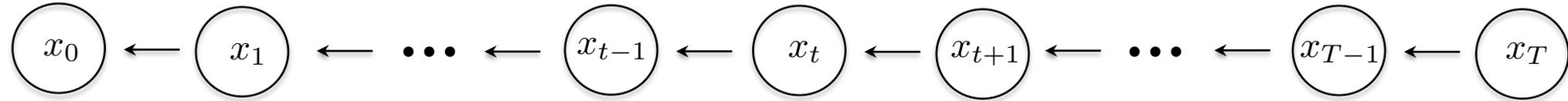
x_n

Outline

- Introduction
- **Vanilla diffusion models (DDPM)**
- Lower bound on data distribution
- Noising process
- Denoising process
- Learning to denoise
- Conclusion

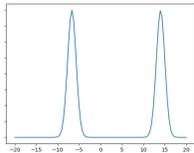
Generative process

- DDPM (Ho-Jain-Abbeel 2020) : Produce new data using a sequence of denoising steps.
Denoising Diffusion Probabilistic Models



$x_0 \sim p(x)$

Probability of training/real data (what we want to estimate)



$x_{t-1} \sim p(x_{t-1}|x_t) = \mathcal{N}(\mu_{\text{denoise}}(x_t, n_t), \sigma_{\text{denoise}}^2(t)\mathbf{I})$
 $x_{t-1} = a_t^x x_t - b_t^n n_t + \sigma_{\text{denoise}}^2(t)z, z \sim \mathcal{N}(0, \mathbf{I})$

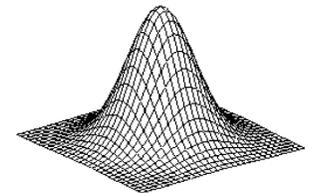
x_{t-1} is produced with a linear combination of the current denoised data x_t , an estimation of the noise n_t and a pure random noise z .

Starting from pure noise $x = x_T$, gradually remove noise to generate intermediate states x_{T-1}, x_{T-2}, \dots until reaching a clean data $x = x_0$, which belongs to the training distribution $p(x)$.

Backward process

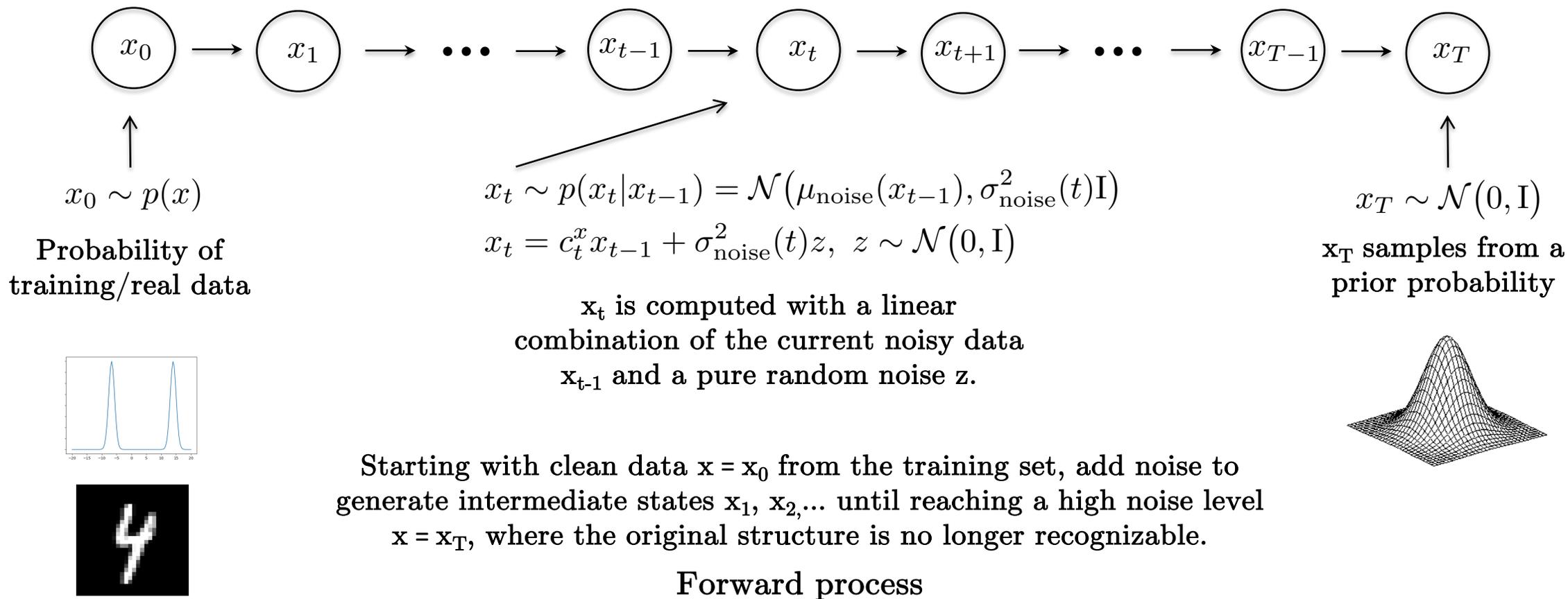
$x_T \sim \mathcal{N}(0, \mathbf{I})$

x_T is sampled from a prior (unconditional) probability.



Forward process

- The denoising steps are learned from the forward pass, which consists in adding noise to the original data :



Outline

- Introduction
- Vanilla diffusion models (DDPM)
- **Lower bound on data distribution**
- Noising process
- Denoising process
- Learning to denoise
- Conclusion

Data distribution

- The goal is to estimate the prior data distribution $p(x)$ from the sequence of noisy steps (forward process) and denoised steps (backward process).
- We start by expressing the data distribution w.r.t. the intermediate states x_1, x_2, \dots :

$$\begin{aligned}\log p(x) &= \log p(x_0), \text{ with } x = x_{t=0} \\ &= \log \int p(x_0, x_1, \dots, x_T) dx_1 dx_2 \dots dx_T, \text{ marginal integration} \\ &= \log \int p(x_{0:T}) dx_{1:T}, \text{ with the sequence notation } x_{0:T} = x_0, x_1, \dots, x_T \\ &= \log \int dx_{1:T} p(x_{0:T}) \frac{p(x_{1:T}|x_0)}{p(x_{1:T}|x_0)}, \text{ where } p(x_{1:T}|x_0) \text{ is the conditional} \\ &\quad \text{probability of having the sequence } x_{1:T} \text{ starting from the original data } x_0.\end{aligned}$$

Data distribution

- Then, we decompose the data distribution in terms of forward and backward probabilities :

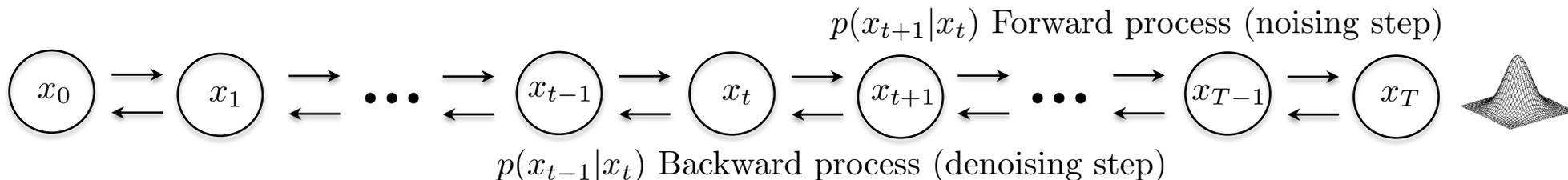
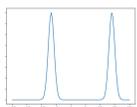
$$\begin{aligned} \log p(x) &= \log \mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\frac{p(x_{0:T})}{p(x_{1:T}|x_0)} \right], \text{ def of expectation: } \int_x p(x) f(x) dx = \mathbb{E}_{x \sim p(x)} [f(x)] \\ &\geq \mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\log \frac{p(x_{0:T})}{p(x_{1:T}|x_0)} \right], \text{ Jensen's inequality: } \varphi(\mathbb{E}[f(x)]) \geq \mathbb{E}[\varphi(f(x))], \varphi \text{ concave} \end{aligned}$$

with chain rule, we have

$$p(x_{0:T}) = p(x_T) p(x_0|x_1) \prod_{t=2}^T p(x_{t-1}|x_t), \text{ given the backward chain process}$$

$$p(x_{1:T}|x_0) = p(x_1|x_0) \prod_{t=2}^T p(x_t|x_{t-1}, x_0), \text{ given the forward chain process}$$

$$\begin{aligned} &\geq \mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\log \frac{p(x_T)p(x_0|x_1)}{p(x_1|x_0)} \right] + \mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\log \prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_t|x_{t-1}, x_0)} \right] \\ &\quad \text{(First term)} \qquad \qquad \qquad \text{(Second term)} \end{aligned}$$



Lower bound

- Let us develop the second term of the lower bound :

$$\begin{aligned} \mathbb{E} \left[\log \prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_t|x_{t-1}, x_0)} \right] &= \mathbb{E} \left[\log \prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{\frac{p(x_{t-1}|x_t, x_0) p(x_t|x_0)}{p(x_{t-1}|x_0)}} \right], \text{ Bayes theorem: } p(a|b)p(b) = p(b|a)p(a) \\ & \qquad \qquad \qquad \text{Here: } p(a|b, x_0)p(b|x_0) = p(b|a, x_0)p(a|x_0) \\ &= \mathbb{E} \left[\log \left(\prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)} \prod_{t=2}^T \frac{p(x_{t-1}|x_0)}{p(x_t|x_0)} \right) \right] \\ &= \mathbb{E} \left[\log \left(\prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)} \frac{p(x_1|x_0)}{p(x_T|x_0)} \right) \right], \text{ terms cancel by recursion} \\ &= \mathbb{E} \left[\underbrace{\log \prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)}} + \mathbb{E} \left[\log \frac{p(x_1|x_0)}{p(x_T|x_0)} \right] \right], \text{ as } \log(ab) = \log a + \log b \end{aligned}$$

A key objective is to reformulate the denoising process $p(x_{t-1}|x_t, x_0)$ to be independent of the clean data x_0 with $p(x_{t-1}|x_t)$ in order to generate new data where x_0 is unknown.

Lower bound

- Focusing on Term 1 and Term 2 :

$$\mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\log p(x_0|x_1) \right] \quad (\text{Term 1})$$

$$= \mathbb{E}_{x_1 \sim p(x_1|x_0)} \left[\log p(x_0|x_1) \right], \text{ function is independent of } x_{2:T}$$

$$\mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\log \frac{p(x_T)}{p(x_T|x_0)} \right] \quad (\text{Term 2})$$

$$= \mathbb{E}_{x_T \sim p(x_T|x_0)} \left[\log \frac{p(x_T)}{p(x_T|x_0)} \right], \text{ function is independent of } x_{1:T-1}$$

$$= -D_{\text{KL}}(p(x_T|x_0), p(x_T)), \text{ definition of Kullback-Leibler divergence}$$

$$D_{\text{KL}}(p_1, p_2) = - \int_x p_1(x) \log \frac{p_2(x)}{p_1(x)} dx = -\mathbb{E}_{x \sim p_1(x)} \left[\log \frac{p_2(x)}{p_1(x)} \right]$$

Lower bound

- Expanding Term 3 :

$$\begin{aligned} & \mathbb{E}_{x_{1:T} \sim p(x_{1:T}|x_0)} \left[\log \prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)} \right] \\ &= \mathbb{E}_{x_{t-1}, x_t \sim p(x_{t-1}, x_t|x_0)} \left[\log \prod_{t=2}^T \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)} \right], \text{ function only depends of } x_{t-1}, x_t \\ &= \sum_{t=2}^T \mathbb{E}_{x_{t-1}, x_t \sim p(x_{t-1}, x_t|x_0)} \left[\log \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)} \right], \text{ as } \log \prod_{t=2}^T = \sum_{t=2}^T \log \\ &= \sum_{t=2}^T \int dx_{t-1} dx_t p(x_{t-1}, x_t|x_0) \log \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)}, \text{ def of expectation} \\ &= \sum_{t=2}^T \int dx_{t-1} dx_t p(x_{t-1}|x_t, x_0) p(x_t|x_0) \log \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)}, \text{ Bayes: } p(a, b|x_0) = p(a|b, x_0)p(b|x_0) \\ &=, - \sum_{t=2}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) \text{ definition of Kullback-Leibler and expectation} \\ & \quad \text{The KL distance is computed in the next slide.} \\ & \text{with } D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) = \int dx_{t-1} p(x_{t-1}|x_t, x_0) \log \frac{p(x_{t-1}|x_t)}{p(x_{t-1}|x_t, x_0)} \\ & \text{and } \mathbb{E}_{x_t \sim p(x_t|x_0)} D_{\text{KL}} = \int dx_t p(x_t|x_0) D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) \end{aligned}$$

Denoising probabilities

- The backward distributions are defined as Gaussian distributions :

Given that

$p(x_{t-1}|x_t, x_0) = \mathcal{N}(\mu_{\text{denoise}_1}(x_t, x_0), \sigma^2(t)\mathbf{I})$, backward pass (denoising process) conditioned on x_t, x_0

$p(x_{t-1}|x_t) = \mathcal{N}(\mu_{\text{denoise}_2}(x_t), \sigma^2(t)\mathbf{I})$, backward pass (denoising process) conditioned on x_t

then the KL distance between these two distributions is

$$D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) = \frac{1}{2\sigma^2(t)} \|\mu_{\text{denoise}_1}(x_t, x_0) - \mu_{\text{denoise}_2}(x_t)\|_2^2$$

$$\text{since } D_{\text{KL}}(p_1, p_2) = \frac{1}{2} \left[\log \left(\frac{\sigma_2^2}{\sigma_1^2} \right) + \frac{\sigma_1^2}{\sigma_2^2} - 1 + \frac{\|\mu_2 - \mu_1\|_2^2}{\sigma_2^2} \right],$$

$$\text{for } p_1 = \mathcal{N}(\mu_1, \sigma_1^2\mathbf{I}), p_2 = \mathcal{N}(\mu_2, \sigma_2^2\mathbf{I}),$$

which reduces to

$$D_{\text{KL}}(p_1, p_2) = \frac{\|\mu_2 - \mu_1\|_2^2}{2\sigma^2} \text{ for two Gaussians with same covariance.}$$

Lower bound

- At this stage, the lower bound estimate of the prior data distribution $p(\mathbf{x})$ is :

$$\begin{aligned} \log p(x) \geq & \underbrace{\mathbb{E}_{x_1 \sim p(x_1|x_0)} \left[\log p(x_0|x_1) \right]}_{\text{(Term 1)}} - \underbrace{D_{\text{KL}}(p(x_T|x_0), p(x_T))}_{\text{(Term 2)}} \\ & - \sum_{t=2}^T \underbrace{\mathbb{E}_{x_t \sim p(x_t|x_0)} D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t))}_{\text{(Term 3)}} \\ & \qquad \qquad \qquad \frac{1}{2\sigma^2(t)} \left\| \mu_{\text{denoise}_1}(x_t, x_0) - \mu_{\text{denoise}_2}(x_t) \right\|_2^2 \end{aligned}$$

with

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}(\mu_{\text{denoise}_1}(x_t, x_0), \sigma^2(t)\mathbf{I}), \text{ backward pass (denoising process)}$$

$$p(x_{t-1}|x_t) = \mathcal{N}(\mu_{\text{denoise}_2}(x_t), \sigma^2(t)\mathbf{I}), \text{ backward pass (denoising process)}$$

Outline

- Introduction
- Vanilla diffusion models (DDPM)
- Lower bound on data distribution
- **Noising process**
- Denoising process
- Learning to denoise
- Conclusion

Noising process

- We design the noising probability (forward process) to be :

$$p(x_t|x_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}),$$

where parameters α_t are the schedulers that control the noise level (explained later).

Using the re-parameterization trick, a sample from $p(x_t|x_{t-1})$ can be expressed as:

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

Noising process

- Let us prove that

$$p(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

$$\text{with the forward probability : } p(x_t|x_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})$$

By recursion, we have

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1}, \epsilon_{t-1} \sim \mathcal{N}(0, \mathbf{I}), \text{ w/ re-parameterization trick}$$

$$x_t = \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2}) + \sqrt{1 - \alpha_t}\epsilon_{t-1}, \epsilon_t, \epsilon_{t-1} \sim \mathcal{N}(0, \mathbf{I})$$

$$x_t = \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \underbrace{\sqrt{\alpha_t}\sqrt{1 - \alpha_{t-1}}\epsilon_{t-2} + \sqrt{1 - \alpha_t}\epsilon_{t-1}}_z$$

Mean of z is zero and the covariance is defined as:

$$\mathbb{E}[zz^T] = [(\sqrt{\alpha_t}\sqrt{1 - \alpha_{t-1}})^2 + (\sqrt{1 - \alpha_t})^2]\mathbf{I} = [1 - \alpha_t\alpha_{t-1}]\mathbf{I}$$

$$x_t = \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\epsilon_{t-2}, \epsilon_{t-2} \sim \mathcal{N}(0, \mathbf{I})$$

After multiple recursions, we finally have

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I}), \text{ with } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

Noising process

- A few observations on the noising process

The noising process $p(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$ with sample

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$$

has an analytical form that only depends on the original clean data x_0 .

This means that no neural network is required to compute a noisy version of x_0 (fast process).

For $t = 0$, $x_{t=0} = x_0$ (original clean image) as $\bar{\alpha}_{t=0} = 1$ and

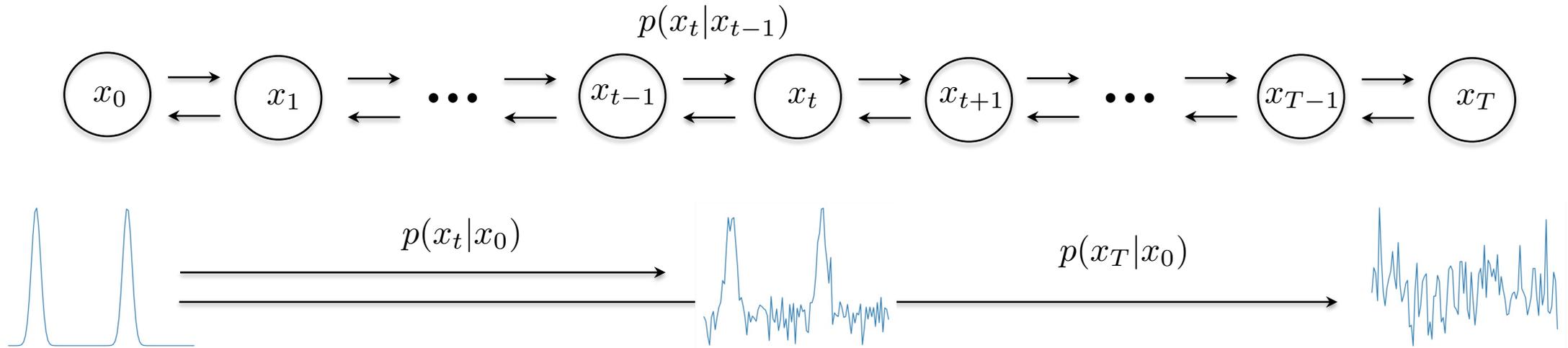
for $t = T$, $x_T = \epsilon_T, \epsilon_T \sim \mathcal{N}(0, \mathbf{I})$ (Gaussian distribution) as $\bar{\alpha}_{t=T} = 0$.

Parameters α_t are called schedulers to transition smoothly from image to pure noise s.a.

$$\alpha_t = \{1, 1 - \delta, 1 - 2\delta, \dots, 1 - t\delta, \dots, 0\}, \text{ where } \delta = \frac{T}{N_t} \text{ is the time step and } N_t \text{ the number of steps.}$$

Noising process

- Illustration :



$$p(x_t|x_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})$$

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

$$p(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I}), \text{ with } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

closed-form solution (no learning is required)

Outline

- Introduction
- Vanilla diffusion models (DDPM)
- Lower bound on data distribution
- Noising process
- **Denoising process**
- Learning to denoise
- Conclusion

Denoising process

- Returning to the backward distribution :

$$\begin{aligned} p(x_{t-1}|x_t, x_0) &= \mathcal{N}(\mu_{\text{denoise}_1}(x_t, x_0), \sigma^2(t)\mathbf{I}) \\ &= \frac{p(x_t|x_{t-1}, x_0)p(x_{t-1}|x_0)}{p(x_t|x_0)}, \text{ Bayes theorem} \end{aligned}$$

with

$$p(x_t|x_{t-1}, x_0) = p(x_t|x_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})$$

$$p(x_{t-1}|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})$$

$$p(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

The backward probability $p(x_{t-1}|x_t, x_0)$ is defined as the product of three Gaussians.

Gaussian product

- The product of the three Gaussians can be written as :

$$\begin{aligned} p(x_{t-1}|x_t, x_0) &\propto \exp\left(-\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{2(1 - \alpha_t)} - \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{2(1 - \bar{\alpha}_{t-1})} + \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{2(1 - \bar{\alpha}_t)}\right) \\ &\propto \exp\left(-\frac{(x_t - \mu_{\text{denoise}_1}(x_t, x_0))^2}{2\sigma^2(t)}\right) \end{aligned}$$

Solving the above quadratic equation provides the mean and the covariance of $p(x_{t-1}|x_t, x_0)$:

$$\begin{aligned} \mu_{\text{denoise}_1}(x_t, x_0) &= \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t}x_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t}x_0 \\ \sigma^2(t) &= \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \end{aligned}$$

Denoising process

- Finally, we arbitrarily design the one-step backward probability $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ by analogy to the mean and the covariance of $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$:

$$p(x_{t-1}|x_t) = \mathcal{N}(\mu_{\text{denoise}_2}(x_t), \sigma^2(t)\mathbf{I})$$

with

$$\mu_{\text{denoise}_2}(x_t) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t}x_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t}\tilde{x}_0(x_t)$$

where $\tilde{x}_0(x_t)$ is the approximate denoised data x_0 from x_t .

$$\text{Recall that } p(x_{t-1}|x_t, x_0) = \mathcal{N}(\mu_{\text{denoise}_1}(x_t, x_0), \sigma^2(t)\mathbf{I})$$

with

$$\mu_{\text{denoise}_1}(x_t, x_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t}x_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t}x_0$$

$$\sigma^2(t) = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}$$

Denoising loss

- Coming back to the KL term :

$$D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) = \frac{1}{\sigma^2(t)} \left\| \mu_{\text{denoise}_1}(x_t, x_0) - \mu_{\text{denoise}_2}(x_t) \right\|_2^2$$

$$\text{with } \mu_{\text{denoise}_1}(x_t, x_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} x_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} x_0$$

$$\text{and } \mu_{\text{denoise}_2}(x_t) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} x_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \tilde{x}_0(x_t)$$

We have

$$D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) = \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \left\| \tilde{x}_0(x_t) - x_0 \right\|_2^2$$

The last equation makes it clear that the backward pass is a denoising process!

When the square error term is minimized then function $\tilde{x}_0(x_t)$ has denoised x_t to be as close as possible to the original clean data x_0 .

Back to maximizing the lower bound

- Going back to the lower bound estimate of the prior data distribution $p(\mathbf{x})$:

$$\begin{aligned} \log p(x) \geq & \mathbb{E}_{x_1 \sim p(x_1|x_0)} \left[\log p(x_0|x_1) \right] - D_{\text{KL}}(p(x_T|x_0), p(x_T)) \\ & \text{(Term 1)} \qquad \qquad \qquad \text{(Term 2)} \\ & - \sum_{t=2}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\tilde{x}_0(x_t) - x_0\|_2^2 \\ & \text{(Term 3)} \end{aligned}$$

with

$$p(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, \sigma^2(t)\mathbf{I})$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$$

and specifically for $t = 1$:

$$x_1 \sim p(x_1|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_1}x_0, \sigma^2(1)\mathbf{I})$$

$$x_1 = \sqrt{\bar{\alpha}_1}x_0 + \sqrt{1 - \bar{\alpha}_1}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$$

Maximizing the lower bound

- Estimating Term 1 :

$$\mathbb{E}_{x_1 \sim p(x_1|x_0)} \left[\log p(x_0|x_1) \right]$$

with $\log p(x_0|x_1) = \log \mathcal{N}(\mu_{\text{denoise}_2}(x_1), \sigma^2(1)\mathbf{I})$

$$\begin{aligned} \text{and } \mu_{\text{denoise}_2}(x_1) &= \frac{(1 - \bar{\alpha}_0)\sqrt{\alpha_1}}{1 - \bar{\alpha}_1} x_1 + \frac{(1 - \alpha_1)\sqrt{\bar{\alpha}_0}}{1 - \bar{\alpha}_1} \tilde{x}_0(x_1) \\ &= \tilde{x}_0(x_1) \text{ as } \alpha_0 = \bar{\alpha}_0 = 1 \text{ and } \alpha_1 = \bar{\alpha}_1 = \delta \end{aligned}$$

which provides

$$\log p(x_0|x_1) \propto \text{ct} - \frac{1}{2\sigma^2(1)} \frac{(1 - \alpha_1)^2 \bar{\alpha}_0}{(1 - \bar{\alpha}_1)^2} \|\tilde{x}_0(x_1) - x_0\|_2^2$$

$$\text{and } - \mathbb{E}_{x_1 \sim p(x_1|x_0)} \frac{1}{2\sigma^2(1)} \frac{(1 - \alpha_1)^2 \bar{\alpha}_0}{(1 - \bar{\alpha}_1)^2} \|\tilde{x}_0(x_1) - x_0\|_2^2 + \text{ct}$$

which is equivalent to Term 3 when $t = 1$ (ct can be ignored when optimizing) :

$$- \sum_{t=2}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\tilde{x}_0(x_t) - x_0\|_2^2$$

Maximizing the lower bound

- Combining Term 1 and Term 3, we have :

$$\log p(x) \geq -D_{\text{KL}}(p(x_T|x_0), p(x_T)) - \sum_{t=1}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\tilde{x}_0(x_t) - x_0\|_2^2$$

with

$$p(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, \sigma^2(t)\mathbf{I})$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$$

Additionally, for Term 2 :

$$p(x_T|x_0) = p(x_T) = \mathcal{N}(0, \mathbf{I}) \quad \forall x_0, \text{ which implies}$$

$$D_{\text{KL}}(p(x_T|x_0), p(x_T)) = 0$$

Finally, we have the lower bound :

$$\log p(x) \geq - \sum_{t=1}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\tilde{x}_0(x_t) - x_0\|_2^2$$

Outline

- Introduction
- Vanilla diffusion models (DDPM)
- Lower bound on data distribution
- Noising process
- Denoising process
- **Learning to denoise**
- Conclusion

MSE loss

- Learning to denoise images :

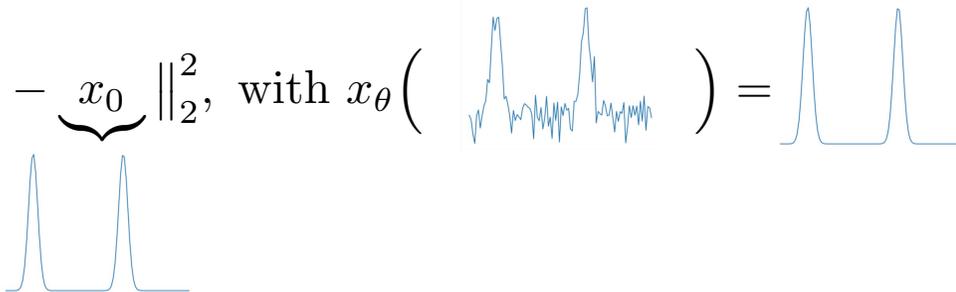
Recall that the noisy process $p(x_t|x_0)$ has a closed-form solution and hence does not require any learning.

However, the denoising function $\tilde{x}_0(\cdot)$ is unknown and needs to be learned with a network $x_\theta(\cdot) = \tilde{x}_0(\cdot)$, where θ are the learnable parameters.

The lower bound function :

$$- \sum_{t=1}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|x_\theta(x_t) - x_0\|_2^2$$

is simplified by dropping the weight coefficients, and the loss to optimize is :

$$L(\theta) = \sum_{t=1}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \|x_\theta(x_t) - \underbrace{x_0}_{\text{clean image}}\|_2^2, \text{ with } x_\theta \left(\text{noisy image} \right) = \text{clean image}$$


and $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$.

Noise prediction

- Instead of predicting the denoised image, an equivalent approach is to predict the noise added to the clean image :

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$$

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_0), \text{ clean image}$$

Substituting x_0 into

$$\mu_{\text{denoise}_1}(x_t, x_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t}x_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t}x_0$$

We have

$$\mu_{\text{denoise}_1}(x_t, x_0) = \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_0$$

As previously, we opt for the simplest design (by analogy) for the mean μ_{denoise_2} :

$$\mu_{\text{denoise}_2}(x_t) = \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_\theta(x_t)$$

where $\epsilon_\theta(x_t)$ is the approximated noise ϵ_0 added to x_0 to produce the noisy image x_t .

Denoising loss

- As previously, considering the KL term :

$$D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) = \frac{1}{2\sigma^2(t)} \left\| \mu_{\text{denoise}_1}(x_t, x_0) - \mu_{\text{denoise}_2}(x_t) \right\|_2^2$$

$$\text{with } \mu_{\text{denoise}_1}(x_t, x_0) = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0$$

$$\text{and } \mu_{\text{denoise}_2}(x_t) = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_\theta(x_t)$$

We now have

$$D_{\text{KL}}(p(x_{t-1}|x_t, x_0), p(x_{t-1}|x_t)) = \frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \left\| \epsilon_\theta(x_t) - \epsilon_0 \right\|_2^2$$

MSE loss

- Learning to predict noise :

The lower bound function becomes

$$- \sum_{t=1}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \left[\frac{1}{2\sigma^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \|\epsilon_\theta(x_t) - \epsilon_0\|_2^2 \right]$$

And the final loss is (dropping the weight coefficients) :

$$L(\theta) = \sum_{t=1}^T \mathbb{E}_{x_t \sim p(x_t|x_0)} \left[\|\epsilon_\theta(x_t) - \underbrace{\epsilon_0}_{\text{noise}}\|_2^2 \right], \text{ with } \epsilon_\theta \left(\begin{array}{c} \text{clean data} \\ \text{+ noise} \end{array} \right) = \text{noise}$$

When the square error term is minimized then function $\epsilon_\theta(x_t)$ predicts the noise ϵ_0 that was added to the original clean data x_0 .

Training steps

- Training a diffusion model to denoise images :

Sample a batch of training data.

Draw random time steps for the batch :

$$t \sim \text{Uniform}([1, \dots, T])$$

Produce noisy samples :

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0, \epsilon_0 \sim \mathcal{N}(0, \mathbf{I})$$

Compute the MSE loss :

$$L(\theta) = \sum_{\text{sampled } t} \|\varepsilon_\theta(x_t) - \varepsilon_0\|_2^2, \text{ with the network}$$

$$\varepsilon_\theta(x) = \text{Transformer/UNet}_\theta(x) \in \mathbb{R}^{b \times d_x}, x \in \mathbb{R}^{b \times d_x}$$

Backpropagation :

Compute gradient of the loss and update net parameters θ .

Generative steps

- Inference : Generate new data

Start with $x_T \sim \mathcal{N}(0, \mathbf{I})$

Compute the sequence $x_{T-1}, x_{T-2}, \dots, x_0$ by sampling at $t = T, T - 1, \dots, 1$

$$x_{t-1} \sim p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_{\text{denoise}_2}(x_t), \sigma^2(t)\mathbf{I})$$

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_\theta(x_t) + \sigma^2(t)z, \quad z \sim \mathcal{N}(0, \mathbf{I})$$

Lab 1 : DDPM for mixture of Gaussians

- DDPM with Transformers for artificial 1D dataset.

```

Jupyter dm_gmm_solution Last Checkpoint: 2 minutes ago
File Edit View Run Kernel Settings Help Not Trusted
JupyterLab Python 3 (ipykernel)

Lab 01 : Diffusion Model (DM) for Mixture of Gaussians -- solution

[1]: # For Google Colaboratory
import sys, os
if 'google.colab' in sys.modules:
    # mount google drive
    from google.colab import drive
    drive.mount('/content/gdrive')
    path_to_file = '/content/gdrive/My Drive/CS242_2025_codes/labs_lecture08/lab03_dm_gmm'
    print(path_to_file)
    # change current path to the folder containing "path_to_file"
    os.chdir(path_to_file)
    !pwd

[3]: # Libraries
import torch
import torch.nn as nn
import torch.optim as optim
import time
# Import utils
import matplotlib.pyplot as plt
import logging
logging.getLogger().setLevel(logging.CRITICAL) # remove warnings
import os, datetime

# PyTorch version and GPU
print(torch.__version__)
if torch.cuda.is_available():
    print(torch.cuda.get_device_name(0))
    device = torch.device("cuda:0") # use GPU
else:
    device = torch.device("cpu")
print(device)

2.2.2
NVIDIA RTX A5000
cuda:0

Create artificial dataset of mixture of Gaussians
    
```

Question 1: Implement the DDPM architecture

Step 1: Define the weights α_t and $\bar{\alpha}_t$

- Their lengths are the same as the number of timesteps T .

Step 2: Code the forward diffusion process

- Jump from x_0 to x_t in one step.

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon_0, \text{ where } \epsilon_0 \sim \mathcal{N}(0, I)$$

Step 3: Implement the backward denoising process

Given the current step sample x_t and t , predict the noise ϵ with a transformer network.

Step 4: Code the generation process

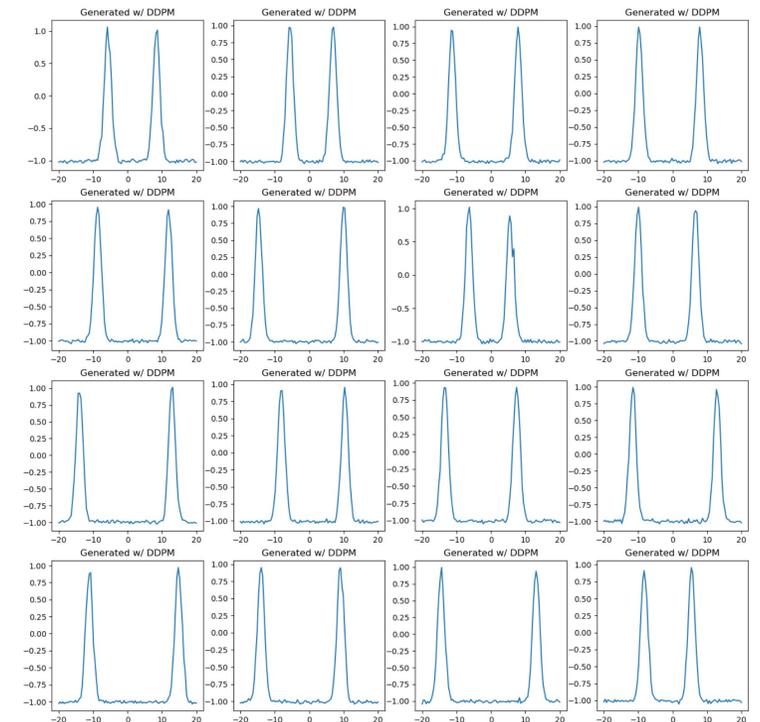
- Randomly sample from $x_T \sim \mathcal{N}(0, I)$.
- Generate x_{t-1} given x_t following the distribution $p(x_{t-1}|x_t, x_0) = \mathcal{N}(\mu_t, \sigma_t^2)$, where:

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

$$\sigma_t^2 = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}$$

$$\text{So, } x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t z, \text{ where } z \sim \mathcal{N}(0, I).$$

We recurrently calculate x_{t-1} until x_0 .



Lab 2 : DDPM for images

- DDPM with CNNs (UNet) for 2D MNIST images.

```

jupyter dm_image_solution Last Checkpoint: 2 minutes ago
File Edit View Run Kernel Settings Help Not Trusted
Python 3 (ipykernel)

Lab 02 : Diffusion Model (DDPM) for MNIST Images -- solution

[1]: # For Google Colaboratory
import sys, os
if 'google.colab' in sys.modules:
    # mount google drive
    from google.colab import drive
    drive.mount('/content/gdrive')
    path_to_file = '/content/gdrive/My Drive/CS5242_2025_codes/labs_lecture08/lab04_dm_image'
    print(path_to_file)
    # move to Google Drive directory
    os.chdir(path_to_file)
    !pwd

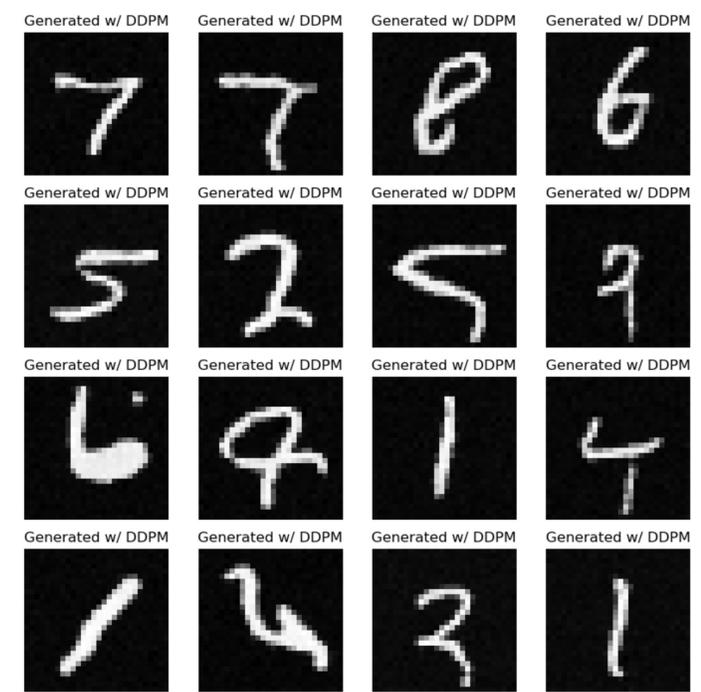
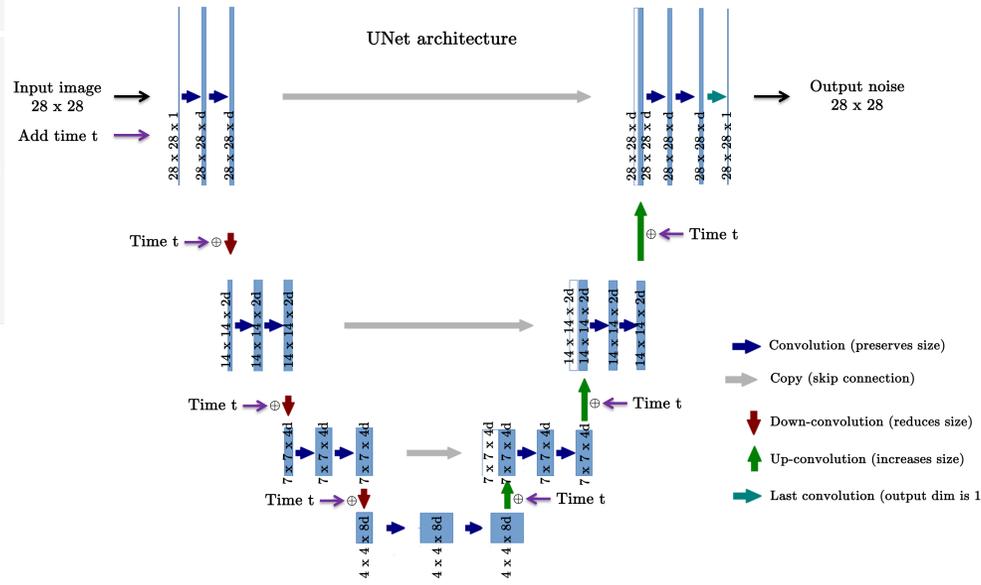
[2]: # Libraries
import torch
import torch.nn as nn
import torch.optim as optim
import time

# import utils
import matplotlib.pyplot as plt
import logging
logging.getLogger().setLevel(logging.CRITICAL) # remove warnings
import os, datetime

# PyTorch version and GPU
print(torch.__version__)
if torch.cuda.is_available():
    print(torch.cuda.get_device_name(0))
    device = torch.device("cuda") # use GPU
else:
    device = torch.device("cpu")
print(device)

2.2.2
NVIDIA RTX A5000
cuda

MNIST dataset
    
```



Outline

- Introduction
- Vanilla diffusion models (DDPM)
- Lower bound on data distribution
- Noising process
- Denoising process
- Learning to denoise
- **Conclusion**

Conclusion

- DDPM (Vanilla Diffusion Model)
 - Surprisingly simple idea!
 - Forward process: Gradually add noise to clean data using a closed-form expression -- this step is fast and analytically tractable.
 - Reverse process: Train a network to denoise the data step-by-step, learning the reverse of the noising process.
 - Generation: Start from pure Gaussian noise and iteratively apply the denoising network to generate new samples.
 - The process is stochastic, so different outputs can be produced from the same initial noise due to the inherent randomness in the model.
 - Most of the math (statistics and algebra) is focused on deriving the correct scaling relationships between the clean image, the denoised prediction, and the added noise.

Conclusion

- Generative models trained on massive datasets are also known as foundation models.
- They have revolutionized text and image processing, powering breakthrough tools like ChatGPT and Stable Diffusion, and driving the rapid growth of the Generative AI (GenAI) industry.
- As these models continue to evolve, especially in their reasoning capabilities, they are becoming increasingly powerful tools -- that will assist us across a wide spectrum of tasks and domains.



Questions?