

# Lecture 01 - Introduction

## ① Main focus:

The principles of diffusion models

- 1. forward process:  
 $\text{data } (x) \rightarrow \text{intermediate distributions} \rightarrow \text{noise } (z)$
- 2. reverse process: (main goal)  
 $\text{noise } (z) \rightarrow \text{intermediate distributions} \rightarrow \text{data } (x)$

Three ways: (regression problem)

- ①. Variational View (VAE)  $\xrightarrow{\text{chp 2}}$  Fokker-Planck Equ.
- ②. Score-based View (EBMs)  $\xrightarrow{\text{chp 3}}$   $\nearrow$
- ③. Flow-based View (Normalizing Flows)  $\xrightarrow{\text{chp 5}}$

- A learned time-dependent velocity field whose flow transports a simple prior to the data.
- Sampling = Solving a differential equation.

Part A: Introduction to DGMs (chp. 1)

Part B: Foundations of DM, (chp. 2-7)

Part C: Sampling of DMs (chp 8-9)

Part D: Learning fast generators (chp. 10-11)

## ②. Deep Generative Modeling

learn a probability distribution  $P_\phi$ .

$$P_\phi \approx P_{\text{data}} \quad \text{via } D_{\text{tr}} := \{x_i\}_{i=1}^N \sim P_{\text{data}}$$

measure:  $D(P_{\text{data}}, P_\phi)$  using  $D_{\text{tr}}$ .

- $D_{\text{tr}} := \{x_1, x_2, \dots, x_N\} : \underbrace{i.i.d}_{\text{indep.}} P_{\text{data}}$

- $P_{\text{data}}$  is intractable

- $P_\phi \approx P_{\text{data}} \quad \text{via}$

$$\phi^* \in \arg \min_{\phi} D(P_{\text{data}}, P_\phi) \quad \text{s.t.}$$

$$P_{\phi^*}(x) \approx P_{\text{data}}(x)$$

choices of  $D$ :

### (1). KL-Divergence

$$D_{\text{KL}}(P_{\text{data}} \| P_\phi) = \int P_{\text{data}}(x) \cdot \log \frac{P_{\text{data}}(x)}{P_\phi(x)} dx$$

$$= \mathbb{E}_{x \sim P_{\text{data}}} [\log P_{\text{data}}(x) - \log P_\phi(x)]$$

- $D_{\text{KL}}(p \| q) \neq D_{\text{KL}}(q, p) \quad \cdot D_{\text{KL}}(p \| q) \geq 0$

$$p = (0.5, 0.5) \quad q = (0.9, 0.1) \quad \cdot D_{\text{KL}}(p \| q) \neq 0$$

$$\begin{cases} D_{\text{KL}}(p \| q) = 0.511 \\ D_{\text{KL}}(q \| p) = 0.367 \end{cases}$$

iff  $p=q$ .

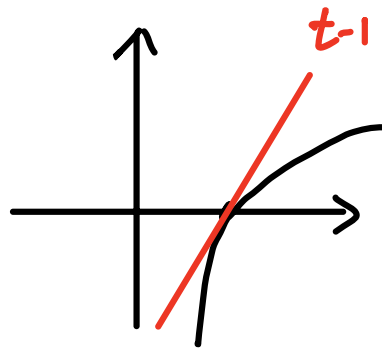
Gibb's Inequality:

$$\begin{cases} D_{KL}(p \parallel q) \geq 0 \\ D_{KL}(p \parallel q) = 0 \text{ iff } p(x) = q(x) \text{ for all } x. \end{cases}$$

proof:  $\forall t > 0, \log t \leq t-1$

$$\log t = t-1 \text{ iff } t=1$$

$$t = \frac{q(x)}{p(x)}$$



$$D_{KL} = \sum_x p(x) \cdot \left( -\log \frac{q(x)}{p(x)} \right)$$

$$-\log t \geq 1-t \Leftrightarrow -\log \frac{q(x)}{p(x)} \geq 1 - \frac{q(x)}{p(x)}$$

$$\Rightarrow p(x) \cdot \left( -\log \frac{q(x)}{p(x)} \right) \geq p(x) - q(x)$$

$$\begin{aligned} \Rightarrow \sum_x p(x) \cdot \left( -\log \frac{q(x)}{p(x)} \right) &\geq \sum_x (p(x) - q(x)) \\ &= \sum_x p(x) - \sum_x q(x) \\ &= 0 \end{aligned}$$

$$\frac{q(x)}{p(x)} = 1 \Leftrightarrow p(x) = q(x) \Leftrightarrow D_{KL}(p \parallel q) = 0. \quad \square$$

(2). Fisher - Divergence:

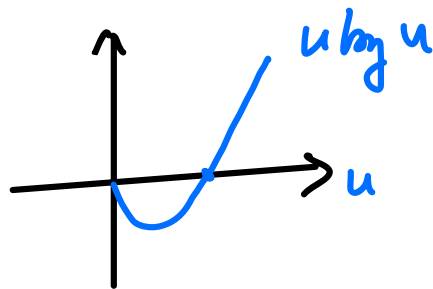
$$D_F(p_{data} \parallel p_\phi) = \mathbb{E}_{x \sim p_{data}} \left[ \underbrace{\|\nabla_x \log p_{data}(x)\|_2^2}_{\text{score fun.}} - \underbrace{\|\nabla_x \log p_\phi(x)\|_2^2}_{\text{score fun.}} \right]$$

• score matching

(3). f-divergence:

$$D_f(p \parallel q) = \int q(x) f\left(\frac{p(x)}{q(x)}\right) dx, \quad f(1) = 0.$$

1.  $f(u) = u \log u$ ,  $D_f = D_{KL}$



2. Jensen-Shannon divergence:

$$D_{JS}(p \parallel q) = \frac{1}{2} D_{KL}(p \parallel m) + \frac{1}{2} D_{KL}(q \parallel m)$$
$$m = \frac{1}{2}(p + q)$$

- $0 \leq D_{JS} \leq \frac{1}{2} \log 2$
- $D_{JS}(p \parallel q) \leq \frac{1}{2} \|p - q\|_1$

3.  $f(u) = \frac{1}{2}|u-1|$ , total variational distance

---

About  $p_\phi(x)$ :

$$\begin{cases} \text{1. } p_\phi(x) \geq 0 \\ \text{2. } \int p_\phi(x) dx = 1. \end{cases}$$

$E_\phi(x)$ : neural network  
 $x \in \mathbb{R}^d$ ,  $E_\phi(x)$

$$\Rightarrow E_\phi(x): \exp(E_\phi(x)) := \tilde{p}_\phi(x)$$

$$Z(\phi) = \int \exp(E_\phi(x')) dx'$$

$$\Rightarrow p_\phi(x) = \frac{\tilde{p}_\phi(x)}{Z(\phi)}$$

③ Examples of DGMs: trade-offs  $\left\{ \begin{array}{l} 1. \text{tractability} \\ 2. \text{expressiveness} \\ 3. \text{training efficiency} \end{array} \right.$

(1). EBM:  $E_{\phi}(x): \mathbb{R}^d \rightarrow \mathbb{R}$

$$p_{\phi}(x) = \frac{\exp(-E_{\phi}(x))}{Z(\phi)}, \quad Z(\phi) = \int \exp(-E_{\phi}(x)) dx$$

Score-based models:

$$\begin{cases} p(x) = e^{f(x)} & \int p(x) dx = 1 \\ q(x) = e^{g(x)} & \int q(x) dx = 1 \end{cases} \quad \text{if } \nabla f = \nabla g \Rightarrow f \equiv g.$$

Proof:  $f(x) = g(x) + c$

$$e^{f(x)} = e^{g(x)+c} = e^c e^{g(x)}$$

$$1 = \int e^{f(x)} dx = \int e^c e^{g(x)} dx = e^c \int e^{g(x)} dx = e^c \cdot 1 = e^c$$

We must have  $c = 0$ .

Score function:  $s(x) = \nabla_x \log p(x)$   $s: \mathbb{R}^d \rightarrow \mathbb{R}^d$ .

$$p(x) = \frac{\tilde{p}(x)}{Z(\phi)}$$

$$\begin{aligned} \nabla_x \log p(x) &= \nabla_x \log \tilde{p}(x) - \underbrace{\nabla_x \log Z}_{=0} \\ &= \nabla_x \log \tilde{p}(x) \end{aligned}$$

## (2) AR-model:

$$P_{\phi}(x) = \prod_{i=1}^D P_{\phi}(x_i | x_{<i})$$

- $x = (x_1, x_2, \dots, x_D)$  sequence data
- generation is slow
- struggle with high-resolution images
- Weak in denoising tasks
- worse for image/video generation,
- Enforce a fixed ordering

But, • works great for text data.

- no need to compute  $Z(\phi)$ .

(3) VAE : learn:  $P_{\phi}$  by decomposing into two parts.

$$\left\{ \begin{array}{l} \text{• encoder (inference network):} \\ \quad q_{\theta}(z|x) \quad x \rightarrow \text{encoder} \rightarrow z \\ \quad \text{(original: } P_{\phi}(z|x)) \quad q_{\theta}(z|x) \\ \text{• Decoder (generator):} \\ \quad P_{\phi}(x|z) \quad z \rightarrow \text{decoder} \rightarrow x' \\ \quad P_{\phi}(x|z) \end{array} \right.$$

⇒ core idea: maximize a lower bound of  $\log P_{\phi}(x)$ .

Evidence Lower Bound (ELBO):

$$P_{\phi}(x) = \int P_{\phi}(x|z) \cdot p(z) dz$$

posterior:  $P_{\phi}(z|x) = \frac{P_{\phi}(x|z) \cdot p(z)}{P_{\phi}(x)} = \frac{P_{\phi}(x|z) \cdot p(z)}{\int P_{\phi}(x|z') dz'}$

intractable

$$q_\theta(z|x) \approx p_\phi(z|x)$$

$$\log p_\phi(x) = \log \int \underline{p_\phi(x, z)} dz$$

$$= \log \int \underbrace{q_\theta(z|x)}_{\text{blue}} \cdot \frac{p_\phi(x, z)}{q_\theta(z|x)} dz$$

$$= \log \mathbb{E}_{z \sim q_\theta(z|x)} \left[ \frac{p_\phi(x, z)}{q_\theta(z|x)} \right]$$

$$\text{// } \log \geq \underbrace{\mathbb{E}_{z \sim q_\theta} \left[ \log \frac{p_\phi(x, z)}{q_\theta(z|x)} \right]}_{\text{ELBO}}$$

$$p_\phi(x, z) = p_\phi(x|z) \cdot p(z)$$

$$= \mathbb{E}_{z \sim q_\theta} \left[ \log \frac{\boxed{p_\phi(x|z)} \cdot \boxed{p(z)}}{\boxed{q_\theta(z|x)}} \right]$$

$$= \underbrace{\mathbb{E}_{z \sim q_\theta} [\log p_\phi(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_\theta(z|x) \parallel p(z))}_{\text{Latent Regularization}}$$

(4) Normalizing flow:

$$\log p_\phi(x) = \log p(z) + \log \left| \det \frac{\partial f_\phi^{-1}(x)}{\partial x} \right|$$

(4) GAN:

$$\left\{ \begin{array}{l} \text{generator: } G_\phi \\ \text{discriminator: } D_\psi \end{array} \right. \quad z \rightarrow G_\phi(z) \rightarrow x'$$

#### ④. Summary:

- Explicit models:  
 $p_\phi(x) \sim$  via tractable or approximately tractable  
ARs, NFs, VAEs, DMs  
define  $p_\phi(x)$  exactly or via a bound.
  - Implicit models:  
Specify a distribution via sampling:  
 $x = G_\phi(z) \quad z \sim p_{\text{prior}}$   
 $p_\phi(x)$ : may not be defined at all.
- NFs, ARs: objective MLE, tractable
  - VAEs, DMs: objective ELBO, Bound / Approximate
  - GANs: objective intractable, not directly modeled



Next topic: variational perspective

(VAE,  $\rightarrow$  DDPM, )

①. VAE

②. DDPM

③. ODE / SDE