

1

# Fundamentals of Generative Diffusion Models

# Form DDPM to Score-based Models

## Fancheng Li (李凡成)

School of Physics and Technology

Wuhan University

Wuhan, China





- O2 A Unified Perspective of Diffusion Models
- O3 Energy-based Models and Guidance
- 04 Stochastic Process and Diffusion Models



- O2 A Unified Perspective of Diffusion Models
- O3 Energy-based Models and Guidance
- •04 Stochastic Process and Diffusion Models

## **Bayes' Rule and Its Validity**





# Our purpose is to learn the posterior probability and likelihood

#### Cromwell's Rule

*"I beseech you, in the bowels of Christ, think it possible that you may be mistaken."* 

I think the moon might be made of cheese

Bernstein-von Mises Theorem

For some case, when n is enough large:

 $||P(\theta|x_1,\ldots,x_n) - \mathcal{N}(\hat{\theta}_n,n^{-1}\mathcal{I}(\theta_0)^{-1})||_{\mathrm{TV}} \xrightarrow{P_{\theta_0}} = 0$ 

# They are the starting point and the end point of Bayesian inference

## **Tweedie's Formula and Evidence Lower Bound**



#### From x to estimate parameters by MSE

 $L = \mathbb{E}[(\hat{\theta}(x) - \theta)^2]$ 

#### from the assumption

$$p(x|\theta) = \mathcal{N}(\theta, \sigma^2)$$
$$p(x) = \int_{-\infty}^{\infty} p(x|\theta) p(\theta) d\theta$$

#### **Tweedie's Estimator**

$$\mathbb{E}[\theta|x] = x + \sigma^2 \frac{\mathrm{d}}{\mathrm{d}x} \log p(x)$$
$$\hat{\theta}^{TE} = x + \sigma^2 \frac{\mathrm{d}}{\mathrm{d}x} \log p(x)$$

 Variational Bayesian inference

 ELBO is a good Loss Function

  $\log p(\boldsymbol{x}) = \int q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) \log p(\boldsymbol{x}) d\boldsymbol{z}$ 
 $= \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \right] - \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{z}|\boldsymbol{x})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \right]$ 
 $= \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \right] + \mathcal{D}_{KL}(q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}|\boldsymbol{x}))$ 
 $= \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \right] + \mathcal{D}_{KL}(q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}|\boldsymbol{x}))$ 
 $= \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p_{\theta}(\boldsymbol{x}|\boldsymbol{z}) \right] - \mathcal{D}_{KL}(q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}|\boldsymbol{x}))$ 

#### $\log p(x)$ is alike energy function in statistical physics and related to score-based or energy-based models

https://en.wikipedia.org/wiki/Bayesian\_statistics



# **• 02 A Unified Perspective of Diffusion Models**

- O3 Energy-based Models and Guidance
- •04 Stochastic Process and Diffusion Models

## **Framework of Diffusion Models**







#### **Forward Process Adds Noise to Images**

$$egin{aligned} & m{x}_t = \sqrt{lpha_t} m{x}_{t-1} + \sqrt{1-lpha_t} m{\epsilon} \ & m{\epsilon} & \sim \mathcal{N}(m{x}_t; 0, m{I}) \end{aligned}$$

# Perturbation from the original image by recursion

$$egin{aligned} & m{x}_t = \sqrt{ar{lpha}_t} m{x}_0 + \sqrt{1 - ar{lpha}_t} m{\epsilon} \ & m{\epsilon} & \sim \mathcal{N}(m{x}_t; \sqrt{ar{lpha}_t} m{x}_0, (1 - ar{lpha}_t)m{I} \end{aligned}$$

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) \coloneqq \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$
$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) \coloneqq \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

**Converting from images to Gaussian noise** 



#### **Reverse Process Denoises the Gaussian Noise**

Generating new images from Gaussian noise by predicting the reverse samples in every time step



## **Evidence Lower Bound (ELBO)**



#### The formula of ELBO

$$egin{aligned} \log p(oldsymbol{x}) &= \log \int rac{p(oldsymbol{x}_{0:T})q(oldsymbol{x}_{1:T}|oldsymbol{x}_{0})}{q(oldsymbol{x}_{1:T}|oldsymbol{x}_{0})} \mathrm{d}oldsymbol{x}_{1:T} \ &= \log \mathbb{E}_{q(oldsymbol{x}_{1:T}|oldsymbol{x}_{0})} \left[rac{p(oldsymbol{x}_{0:T})}{q(oldsymbol{x}_{1:T}|oldsymbol{x}_{0})}
ight] \ &\geq \mathbb{E}_{q(oldsymbol{x}_{1:T}|oldsymbol{x}_{0})} \left[\log rac{p(oldsymbol{x}_{0:T})}{q(oldsymbol{x}_{1:T}|oldsymbol{x}_{0})}
ight] \end{aligned}$$

The model is trained

#### by maximizing ELBO

So we need to further learn about ELBO ELBO can be written as three terms:

$$= \underbrace{\mathbb{E}_{q(\boldsymbol{x}_{1}|\boldsymbol{x}_{0})} \left[\log p_{\theta}(\boldsymbol{x}_{0}|\boldsymbol{x}_{1})\right]}_{\text{ Edam}} \\ - \underbrace{D_{KL}(q(\boldsymbol{x}_{T}|\boldsymbol{x}_{0}) \parallel p(\boldsymbol{x}_{T}))}_{\text{ Ewmann}} \\ - \underbrace{\sum_{t=2}^{T} \mathbb{E}_{q(\boldsymbol{x}_{t},\boldsymbol{x}_{t-1}|\boldsymbol{x}_{0})} \left[D_{KL}(q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_{t},\boldsymbol{x}_{0}) \parallel p_{\theta}(\boldsymbol{x}_{t-1}|\boldsymbol{x}_{t}))\right]}_{\text{ Ewmann}} \\ = \underbrace{\sum_{t=2}^{T} \mathbb{E}_{q(\boldsymbol{x}_{t},\boldsymbol{x}_{t-1}|\boldsymbol{x}_{0})} \left[D_{KL}(q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_{t},\boldsymbol{x}_{0}) \parallel p_{\theta}(\boldsymbol{x}_{t-1}|\boldsymbol{x}_{t}))\right]}_{\text{ Ewmann}}$$

The first two terms are determined by forward process.

So we only need to optimize the last term.

### **Several Equivalent Optimization Function**



#### For KL divergence of Gaussian distribution

$$egin{aligned} &D_{KL}(\mathcal{N}(oldsymbol{x};oldsymbol{\mu_x},oldsymbol{\Sigma_y}) \|\mathcal{N}(oldsymbol{y};oldsymbol{\mu_y},oldsymbol{\Sigma_y})) \ &=&rac{1}{2}iggl[ \lograc{|oldsymbol{\Sigma_y}|}{|oldsymbol{\Sigma_x}|}+tr(oldsymbol{\Sigma_y}^{-1}oldsymbol{\Sigma_x})+(oldsymbol{\mu_y}-oldsymbol{\mu_x})^{\mathrm{T}}oldsymbol{\Sigma_y}^{-1}(oldsymbol{\mu_y}-oldsymbol{\mu_x})-n \ &oldsymbol{\omega_x}\ &&\mbol{arg}\min_{oldsymbol{ heta}}D_{KL}(q(oldsymbol{x}_{t-1}|oldsymbol{x}_{t},oldsymbol{x}_{0})\parallel p_{oldsymbol{ heta}}(oldsymbol{x}_{t-1}|oldsymbol{x}_{t})) \ &oldsymbol{\omega_x}\ &&\mbol{arg}\min_{oldsymbol{ heta}}rac{1}{2\sigma_q^2(t)}iggl[\|oldsymbol{\mu}_{ heta}-oldsymbol{\mu}_{q}\|_2^2iggr] \end{aligned}$$

So we get the first **optimization function** we need to predict the **mean value** by neural networks Predict samples by neural networks

$$oldsymbol{\mu}_{oldsymbol{ heta}}(oldsymbol{x}_t,t) = rac{\sqrt{lpha_t}(1-ar lpha_{t-1})oldsymbol{x}_t + \sqrt{ar lpha_{t-1}}(1-lpha_t)oldsymbol{\hat x}_{oldsymbol{ heta}}(oldsymbol{x}_t,t)}{1-ar lpha_t} \ \|oldsymbol{a}_{q}(oldsymbol{x}_t,t) - oldsymbol{x}_{0})\|_2^2 igg|$$

**Predict noise by neural networks** 

$$oldsymbol{\mu}_{ heta}(oldsymbol{x}_t,t) = rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\hat{\epsilon}}_{oldsymbol{ heta}}(oldsymbol{x}_t,t) 
onumber \ \mathbf{x}_t,t) 
onumber \ \mathbf{x}_t = rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t + rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\hat{\epsilon}}_{oldsymbol{ heta}}(oldsymbol{x}_t,t) 
onumber \ \mathbf{x}_t,t) 
onumber$$

### **Several Equivalent Optimization Function**



#### **Tweedie's Estimator**

$$\mathbb{E}[\theta|x] = x + \sigma^2 \frac{\mathrm{d}}{\mathrm{d}x} \log p(x)$$

$$\mathbb{E}[\boldsymbol{\mu}_{\boldsymbol{x}_t}|\boldsymbol{x}_t] = \boldsymbol{x}_t + (1 - \bar{\alpha}_t)\nabla_{\boldsymbol{x}_t} \log p(\boldsymbol{x}_t)$$

$$\sqrt{\bar{\alpha}_t}\boldsymbol{x}_0 = \boldsymbol{x}_t + (1 - \bar{\alpha}_t)\nabla_{\boldsymbol{x}_t} \log p(\boldsymbol{x}_t)$$

$$\boldsymbol{x}_0 = \frac{\boldsymbol{x}_t + (1 - \bar{\alpha}_t)\nabla_{\boldsymbol{x}_t} \log p(\boldsymbol{x}_t)}{\sqrt{\bar{\alpha}_t}}$$

$$\boldsymbol{\mu}_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\boldsymbol{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}}s_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t)$$

$$\arg\min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} \Big[ \|s_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t) - \nabla \log p(\boldsymbol{x}_t)\|_2^2$$

Relation between noise and score function

$$abla_{oldsymbol{x}_t}\log p(oldsymbol{x}_t) = -rac{1}{\sqrt{1-ar{lpha}_t}}oldsymbol{\epsilon}_0$$

The score function measures how the data is move to maximize, it clear that the opposite direction is noisy and as can be seen from the above formula.





O2 A Unified Persperctive of Diffusion Models

# **O3 Energy-based Models and Guidance**

## •04 Stochastic Process and Diffusion Models

### **Energy-based Models**



#### Why the last formula can lead to energy-based model

 $egin{aligned} p_{ heta}(oldsymbol{x}) &= rac{1}{z_{ heta}} e^{-f_{ heta}(oldsymbol{x})} \ 
abla_{oldsymbol{x}} \log p_{oldsymbol{ heta}}(oldsymbol{x}) &= 
abla_{oldsymbol{x}} \log rac{1}{z_{oldsymbol{ heta}}} e^{-f_{oldsymbol{ heta}}(oldsymbol{x})} \ 
&= 
abla_{oldsymbol{x}} \log rac{1}{z_{oldsymbol{ heta}}} + 
abla_{oldsymbol{x}} \log e^{-f_{oldsymbol{ heta}}(oldsymbol{x})} \ 
&= abla_{oldsymbol{x}} \log rac{1}{z_{oldsymbol{ heta}}} + 
abla_{oldsymbol{x}} \log e^{-f_{oldsymbol{ heta}}(oldsymbol{x})} \ 
&= abla_{oldsymbol{x}} f_{oldsymbol{ heta}}(oldsymbol{x}) \ 
&\approx s_{oldsymbol{ heta}}(oldsymbol{x}) \ 
&\equiv s_{oldsymbol{ heta}}(oldsymbol{ heta}) \ 
&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 
&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 
&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{ heta}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{ heta}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{x}) \ 

&\equiv s_{oldsymbol{ heta}}(oldsymbol{ heta}) \ 

&\equiv s_{oldsymbol{$ 

We can further see the relationship between the diffusion models and statistical physics.

#### Langevin dynamics

$$oldsymbol{x}_{i+1} = oldsymbol{x}_i + c 
abla \log p(oldsymbol{x}_i) + \sqrt{2c}oldsymbol{\epsilon}_i$$

The score function represents a move on the manifold

Noise protection against **local optimality** 



### Guidance



#### **Classifier Guidance**

$$egin{aligned} 
abla \log p(oldsymbol{x}_t|y) &= 
abla \log \left(rac{p(oldsymbol{x}_t)p(y|oldsymbol{x}_t)}{p(y)}
ight) \ &= 
abla \log p(oldsymbol{x}_t) + 
abla \log p(y|oldsymbol{x}_t) - 
abla \log p(y) \ &= \underbrace{
abla \log p(oldsymbol{x}_t)}_{ ext{T}\& ext{H} lpha eta} + \underbrace{
abla \log p(y|oldsymbol{x}_t)}_{ ext{H} rac{1}{2} ext{H} lpha eta} \end{aligned}$$

 $abla \log p(oldsymbol{x}_t|y) = 
abla \log p(oldsymbol{x}_t) + \gamma 
abla \log p(y|oldsymbol{x}_t)$ 

*Classifer guidance need to train two different diffusion models* 

So training is **very expensive.** 

#### **Classifier-Free Guidance**

$$abla \log p(y|oldsymbol{x}_t) = 
abla \log p(oldsymbol{x}_t|y) - 
abla \log p(oldsymbol{x}_t)$$

$$egin{aligned} 
abla \log p(oldsymbol{x}_t|y) &= 
abla \log p(oldsymbol{x}_t) + \gamma(
abla \log p(oldsymbol{x}_t|y) - 
abla \log p(oldsymbol{x}_t)) \ &= \underbrace{\gamma 
abla \log p(oldsymbol{x}_t|y)}_{\Re atheta 
angle + \underbrace{(1 - \gamma) 
abla \log p(oldsymbol{x}_t)}_{ atheta 
angle + rac{(1 - \gamma) 
abla \log p(oldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \log p(oldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \log p(oldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \log p(oldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \log p(oldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}_t)}{ atheta 
angle + rac{(1 - \gamma) 
abla \boldsymbol{x}_t \boldsymbol{x}$$

Classifer-Free guidance can train the diffusion models with guidance with **more convenience.** We can get the output of diffusion models with **more requirements.** 



- O2 A Unified Persperctive of Diffusion Models
- O3 Energy-based Models and Guidance

# O4 Stochastic Process and Diffusion Models

## **Stochastic Differential Equations**





### From SDE to DDPM



It is easy to show that DDPM can be seen as a special case of SDE.

$$egin{aligned} oldsymbol{x}(t+\Delta t) &= \sqrt{1-eta(t+\Delta t)\Delta t}oldsymbol{x}(t) + \sqrt{eta(t+\Delta t)\Delta t}oldsymbol{arepsilon}(t) \ oldsymbol{x}(t+\Delta t) &pprox \left[1-rac{eta(t+\Delta t)\Delta t}{2}
ight]oldsymbol{x}(t) + \sqrt{eta(t+\Delta t)\Delta t}oldsymbol{arepsilon}(t) \ oldsymbol{d}oldsymbol{x}(t) &pprox -rac{eta(t)oldsymbol{x}(t)}{2}oldsymbol{d}t + \sqrt{eta(t)}\sqrt{oldsymbol{d}t}oldsymbol{arepsilon}(t) \ oldsymbol{d}oldsymbol{x} &= -rac{eta(t)oldsymbol{x}(t)}{2}oldsymbol{d}t + \sqrt{eta(t)}oldsymbol{d}oldsymbol{w} \end{aligned}$$

SDE provides a lot of convenience for design and give an understanding of diffusion model of stochastic thermodynamics



Group Report



# Thanks for listening

Fancheng Li (李凡成)