

Generative Model With Variational Inference —VAE and Diffusion Model

Yuzhou Nie(聂宇舟)

School of Statistics, Renmin University of China



中國人民大學
RENMIN UNIVERSITY OF CHINA



Likelihood-based Model

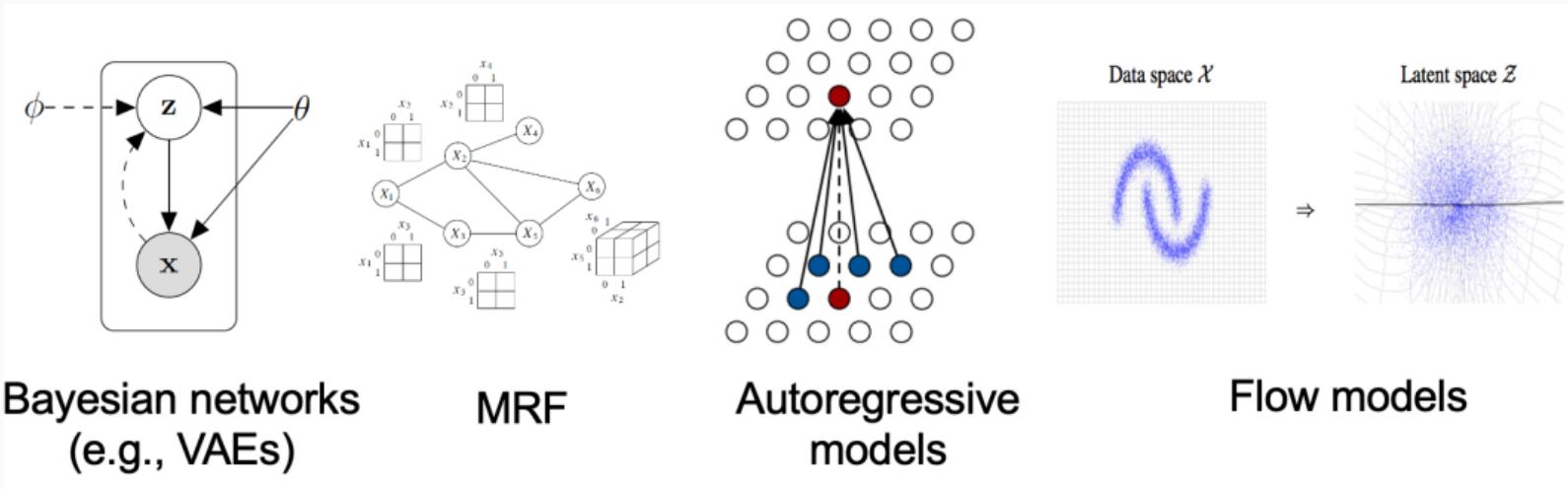


Figure 1: Typical Likelihood-based Models

Likelihood-based Model

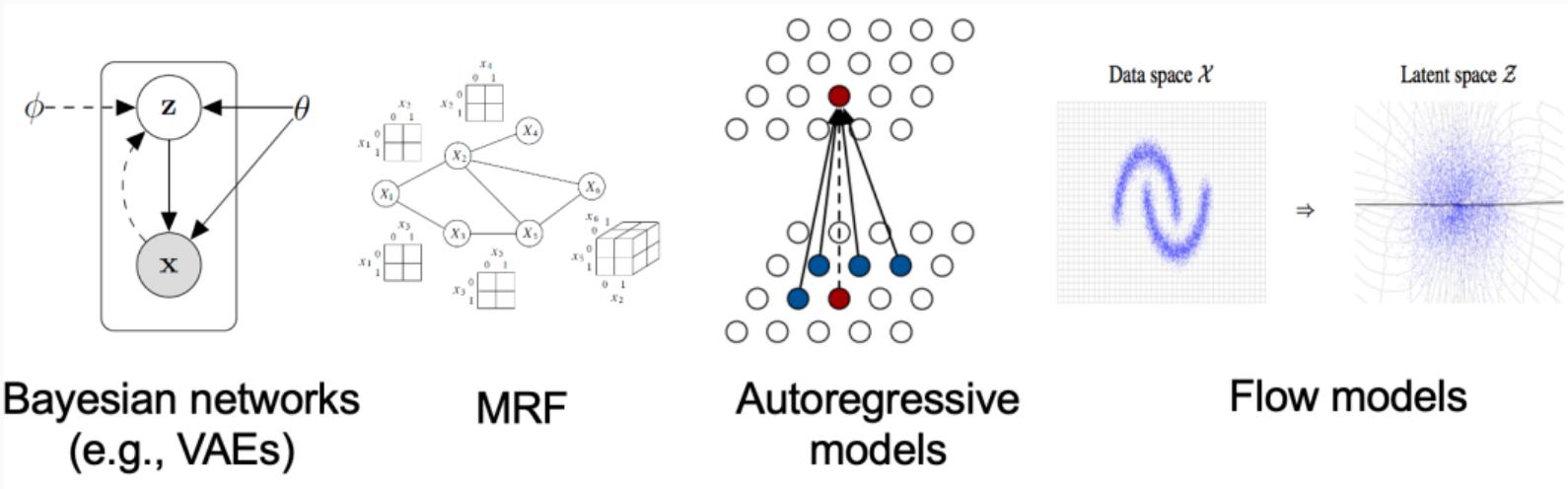


Figure 1: Typical Likelihood-based Models

Variational Inference

Contents

- Variational Inference
- Variational Autoencoder and Expectation Maximum
- Diffusion Model: "Autoregressive" VAE
- Score-based Models Based on SDEs: A General Framework

Some Important Notes

- Skip the details of derivation.

Some Important Notes

- Skip the details of derivation.
- Skip the details of basic proposition and conclusion in statistics.

Some Important Notes

- Skip the details of derivation.
- Skip the details of basic proposition and conclusion in statistics.
- Please remember important formulae.

Some Important Notes

- Skip the details of derivation.
- Skip the details of basic proposition and conclusion in statistics.
- Please remember important formulae.
- $p_\theta(x)$ can be regarded as $p(x | \theta)$ here.

Start From Maximum Likelihood Estimation

Suppose x are observed variables, $p(x)$ is empirical distribution of x , $p_\theta(x)$ is any parameterized distribution of x .

Start From Maximum Likelihood Estimation

Suppose x are observed variables, $p(x)$ is empirical distribution of x , $p_\theta(x)$ is any parameterized distribution of x .

$$p_\theta(x) \rightarrow p(x)$$

Start From Maximum Likelihood Estimation

Suppose x are observed variables, $p(x)$ is empirical distribution of x , $p_\theta(x)$ is any parameterized distribution of x .

$$p_\theta(x) \rightarrow p(x)$$

Maximize log likelihood estimation.

$$\theta = \arg \max_{\theta} \mathbb{E}_{x \sim p(x)} [\log p_\theta(x) dx] = \arg \max_{\theta} \int p(x) \log p_\theta(x) dx \quad (1)$$

Start From Maximum Likelihood Estimation

Suppose x are observed variables, $p(x)$ is empirical distribution of x , $p_\theta(x)$ is any parameterized distribution of x .

$$p_\theta(x) \rightarrow p(x)$$

Maximize log likelihood estimation.

$$\theta = \arg \max_{\theta} \mathbb{E}_{x \sim p(x)} [\log p_\theta(x) dx] = \arg \max_{\theta} \int p(x) \log p_\theta(x) dx \quad (1)$$

Minimize KL divergence $\text{KL}(p(x) \| p_\theta(x))$:

$$\theta = \arg \min_{\theta} \text{KL}(p(x) \| p_\theta(x)) = \arg \min_{\theta} \int p(x) \log \left(\frac{p(x)}{p_\theta(x)} \right) dx \quad (2)$$

Start From Maximum Likelihood Estimation

Suppose x are observed variables, $p(x)$ is empirical distribution of x , $p_\theta(x)$ is any parameterized distribution of x .

$$p_\theta(x) \rightarrow p(x)$$

Maximize log likelihood estimation.

$$\theta = \arg \max_{\theta} \mathbb{E}_{x \sim p(x)} [\log p_\theta(x) dx] = \arg \max_{\theta} \int p(x) \log p_\theta(x) dx \quad (1)$$

Minimize KL divergence $\text{KL}(p(x) \| p_\theta(x))$:

$$\theta = \arg \min_{\theta} \text{KL}(p(x) \| p_\theta(x)) = \arg \min_{\theta} \int p(x) \log \left(\frac{p(x)}{p_\theta(x)} \right) dx \quad (2)$$

Hard to optimize

Variational Inference

Introduce z as unobserved variables, we infer $p(z | x)$

Variational Inference

Introduce z as unobserved variables, we infer $p(z | x)$

$$p(z | x) = \frac{p(x, z)}{\underbrace{\int_z p(x, z) dz}_{\text{intractable}}} \quad (3)$$

Variational Inference

Introduce z as unobserved variables, we infer $p(z | x)$

$$p(z | x) = \frac{p(x, z)}{\underbrace{\int_z p(x, z) dz}_{\text{intractable}}} \quad (3)$$

Possible solution: $q_\phi(z | x) \rightarrow p(z | x)$

Variational Inference

Introduce z as unobserved variables, we infer $p(z | x)$

$$p(z | x) = \frac{p(x, z)}{\underbrace{\int_z p(x, z) dz}_{\text{intractable}}} \quad (3)$$

Possible solution: $q_\phi(z | x) \rightarrow p(z | x)$

Minimize KL divergence:

$$\phi = \arg \min_{\phi} \mathbb{KL}(q_\phi(z | x) \| p(z | x)) \quad (4)$$

Variational Inference

Introduce z as unobserved variables, we infer $p(z | x)$

$$p(z | x) = \frac{p(x, z)}{\underbrace{\int_z p(x, z) dz}_{\text{intractable}}} \quad (3)$$

Possible solution: $q_\phi(z | x) \rightarrow p(z | x)$

Minimize KL divergence:

$$\phi = \arg \min_{\phi} \mathbb{KL}(q_\phi(z | x) \| p(z | x)) \quad (4)$$

$$\begin{aligned} \mathbb{KL}(q_\phi(z | x) \| p(z | x)) &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z | x)}{p(z | x)} \right) \right] \\ &= \log p(x) - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p(x, z)}{q_\phi(z | x)} \right) \right] \end{aligned} \quad (5)$$

Variational Inference

$$\min \mathbb{KL}(q_\phi(z | x) \| p(z | x)) = \max \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p(x,z)}{q_\phi(z | x)} \right) \right]}_{ELBO(q_\phi)} \quad (6)$$

Variational Inference

$$\min \mathbb{KL}(q_\phi(z | x) \| p(z | x)) = \max \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p(x,z)}{q_\phi(z | x)} \right) \right]}_{ELBO(q_\phi)} \quad (6)$$

Evidence Lower Bound

Variational Inference

$$\min \mathbb{KL}(q_\phi(z | x) \| p(z | x)) = \max \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p(x,z)}{q_\phi(z | x)} \right) \right]}_{ELBO(q_\phi)} \quad (6)$$

Evidence Lower Bound

Further derive:

$$\begin{aligned} ELBO(q_\phi) &= \mathbb{E}_{q_\phi} [\log p(z)] + \mathbb{E}_{q_\phi} [\log p(x | z)] - \mathbb{E}_{q_\phi} [\log q_\phi(z | x)] \\ &= \underbrace{\mathbb{E}_{q_\phi} [\log p(x | z)]}_{likelihood} - \underbrace{\mathbb{KL}(q_\phi(z | x) \| p(z))}_{\substack{\text{variational} \\ \text{prior}}} \end{aligned} \quad (7)$$

Variational Inference From Generative Model

$$\log p_{\theta}(x) \rightarrow \log p(x)$$

Variational Inference From Generative Model

$$\log p_\theta(x) \rightarrow \log p(x)$$

suppose $p_\theta(x) = \int p_\theta(x, z) dz$, we have

$$\begin{aligned} \log p_\theta(x) &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{p_\theta(z|x)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \right) \right]}_{ELBO(q_\phi)} + \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right]}_{\text{KL}(q_\phi(z|x) \| p_\theta(z|x))} \end{aligned} \tag{8}$$

Variational Inference From Generative Model

$$\log p_\theta(x) \rightarrow \log p(x)$$

suppose $p_\theta(x) = \int p_\theta(x, z) dz$, we have

$$\begin{aligned}\log p_\theta(x) &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{p_\theta(z|x)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \right) \right]}_{ELBO(q_\phi)} + \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right]}_{\text{KL}(q_\phi(z|x) \| p_\theta(z|x))}\end{aligned}\tag{8}$$

$\log p_\theta(x) \geq ELBO(q_\phi)$, equal when $q_\phi(z|x) = p_\theta(z|x)$

Variational Inference From Generative Model

$$\log p_\theta(x) \rightarrow \log p(x)$$

suppose $p_\theta(x) = \int p_\theta(x, z) dz$, we have

$$\begin{aligned}\log p_\theta(x) &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{p_\theta(z|x)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right] \\ &= \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \right) \right]}_{ELBO(q_\phi)} + \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right]}_{\text{KL}(q_\phi(z|x) \| p_\theta(z|x))}\end{aligned}\tag{8}$$

$\log p_\theta(x) \geq ELBO(q_\phi)$, equal when $q_\phi(z|x) = p_\theta(z|x)$

ELBO is a lower bound for $\log p_\theta(x)$

Variational Inference

For KL divergence:

Variational Inference

For KL divergence:

$$\begin{aligned} \log p_\theta(x) &\geq ELBO(q_\phi) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)})] \\ \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)p(x)})]] \end{aligned} \tag{9}$$

Variational Inference

For KL divergence:

$$\begin{aligned} \log p_\theta(x) &\geq ELBO(q_\phi) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)})] \\ \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)p(x)})]] \end{aligned} \tag{9}$$

According to $q_\phi(z|x) \rightarrow p(z|x)$, we have $p(x,z) = p(z|x)p(x) \approx q_\phi(z|x)p(x)$

Variational Inference

For KL divergence:

$$\begin{aligned} \log p_\theta(x) &\geq ELBO(q_\phi) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)})] \\ \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)p(x)})]] \end{aligned} \quad (9)$$

According to $q_\phi(z|x) \rightarrow p(z|x)$, we have $p(x,z) = p(z|x)p(x) \approx q_\phi(z|x)p(x)$

$$\begin{aligned} \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)p(x)})]] \\ \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x,z \sim p(x,z)} [\log \frac{p_\theta(x,z)}{p(x,z)}] \\ \mathbb{KL}(p(x) \| p_\theta(x)) &\leq \mathbb{KL}(p(x,z) \| p_\theta(x,z)) \end{aligned} \quad (10)$$

Variational Inference

For KL divergence:

$$\begin{aligned} \log p_\theta(x) &\geq ELBO(q_\phi) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)})] \\ \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)p(x)})]] \end{aligned} \quad (9)$$

According to $q_\phi(z|x) \rightarrow p(z|x)$, we have $p(x,z) = p(z|x)p(x) \approx q_\phi(z|x)p(x)$

$$\begin{aligned} \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q_\phi(z|x)} [\log(\frac{p_\theta(x,z)}{q_\phi(z|x)p(x)})]] \\ \mathbb{E}_{x \sim p(x)} [\log \frac{p_\theta(x)}{p(x)}] &\geq \mathbb{E}_{x,z \sim p(x,z)} [\log \frac{p_\theta(x,z)}{p(x,z)}] \\ \mathbb{KL}(p(x) \| p_\theta(x)) &\leq \mathbb{KL}(p(x,z) \| p_\theta(x,z)) \end{aligned} \quad (10)$$

$\mathbb{KL}(p(x,z) \| p_\theta(x,z))$ **is an upper bound for** $\mathbb{KL}(p(x) \| p_\theta(x))$

VAE Framework

VAE Framework

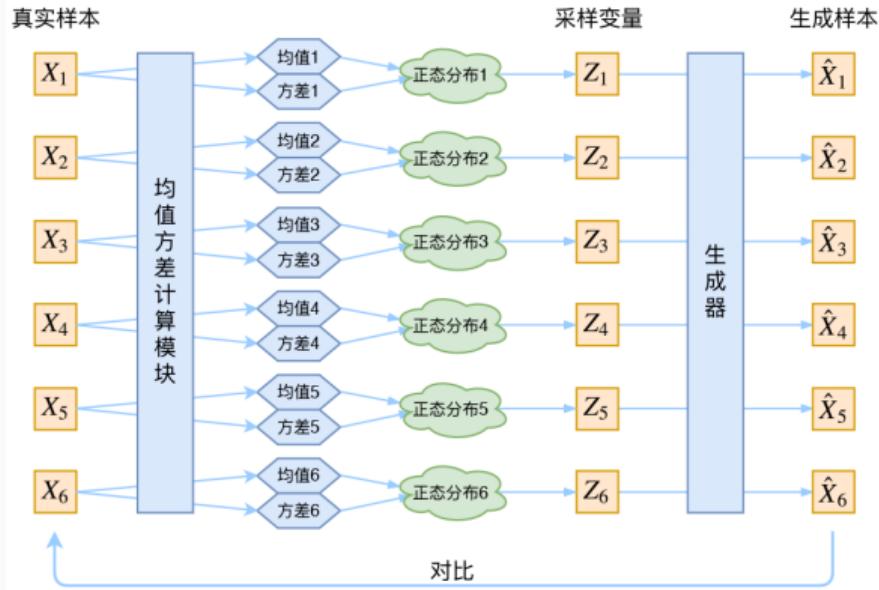


Figure 2: Variational Autoencoder

VAE Framework

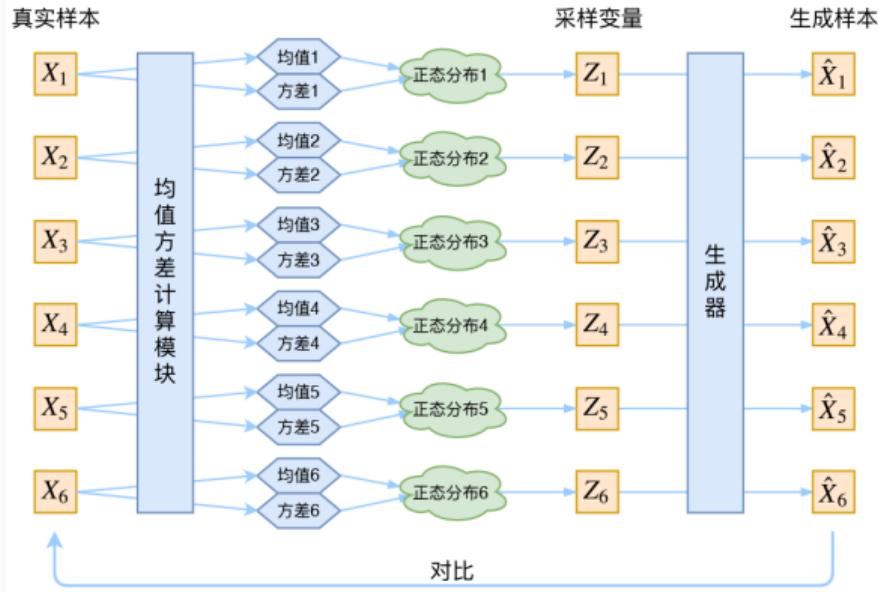


Figure 2: Variational Autoencoder

Objective: $\text{loss} = \text{reconstruction loss}(X, \hat{X}) + \text{regularization}(\mathbb{KL}(N(\mu, \sigma), N(\mathbf{o}, \mathbf{1})))$

VAE

$\mathbb{KL}(p(x, z) \| p_\theta(x, z))$ is an upper bound for $\mathbb{KL}(p(x) \| p_\theta(x))$

VAE

$\mathbb{KL}(p(x, z) \| p_\theta(x, z))$ is an upper bound for $\mathbb{KL}(p(x) \| p_\theta(x))$

We assume

$$p_\theta(x, z) = p_\theta(x | z)q(z), p(x, z) = p(x)q_\phi(z | x)$$

where $p_\theta(x | z)$, $q_\phi(z | x)$ are Gaussian distribution with unknown parameters, $q(z)$ is standard Gaussian distribution.

VAE

$\mathbb{KL}(p(x, z) \| p_\theta(x, z))$ is an upper bound for $\mathbb{KL}(p(x) \| p_\theta(x))$

We assume

$$p_\theta(x, z) = p_\theta(x | z)q(z), p(x, z) = p(x)q_\phi(z | x)$$

where $p_\theta(x | z), q_\phi(z | x)$ are Gaussian distribution with unknown parameters, $q(z)$ is standard Gaussian distribution.

Remark

With $p_\theta(x | z)$ and $q(z)$, the better approximation for $p(z | x)$ should be:

$$p(z | x) = \frac{p_\theta(x | z)q(z)}{p(x)} = \frac{p_\theta(x | z)q(z)}{\underbrace{\int p_\theta(x | z)q(z)dz}_{intractable}} \quad (11)$$

VAE Loss Function

Objective function:

$$\mathbb{KL}(p(x, z) \| p_{\theta}(x, z)) = \iint p(x)q_{\phi}(z | x) \log \frac{p(x)q_{\phi}(z | x)}{p_{\theta}(x | z)q(z)} dx dz \quad (12)$$

VAE Loss Function

Objective function:

$$\mathbb{KL}(p(x, z) \| p_{\theta}(x, z)) = \iint p(x) q_{\phi}(z | x) \log \frac{p(x) q_{\phi}(z | x)}{p_{\theta}(x | z) q(z)} dx dz \quad (12)$$

$\log p(x)$ is irrelevant, we have

$$\mathbb{E}_{x \sim p(x)} \left[- \int q_{\phi}(z | x) \log p_{\theta}(x | z) dz + \mathbb{KL}(q_{\phi}(z | x) \| q(z)) \right] \quad (13)$$

VAE Loss Function

Objective function:

$$\mathbb{KL}(p(x, z) \| p_{\theta}(x, z)) = \iint p(x) q_{\phi}(z | x) \log \frac{p(x) q_{\phi}(z | x)}{p_{\theta}(x | z) q(z)} dx dz \quad (12)$$

$\log p(x)$ is irrelevant, we have

$$\mathbb{E}_{x \sim p(x)} \left[- \int q_{\phi}(z | x) \log p_{\theta}(x | z) dz + \mathbb{KL}(q_{\phi}(z | x) \| q(z)) \right] \quad (13)$$

If we generate one sample for $q_{\phi}(z | x)$ each time

$$\mathbb{E}_{x \sim p(x)} \left[\underbrace{- \log p_{\theta}(x | z)}_{\text{reconstruction}} + \underbrace{\mathbb{KL}(q_{\phi}(z | x) \| q(z))}_{\text{regularization}} \right] \quad (14)$$

EM Algorithm

Assumption for VAE: $p_\theta(x | z)$, $q_\phi(z | x)$ are Gaussian distributions.

EM Algorithm

Assumption for VAE: $p_\theta(x | z)$, $q_\phi(z | x)$ are Gaussian distributions.

Directly optimize $\text{KL}(p(x, z) \| p_\theta(x, z))$

EM Algorithm

Directly optimize $\mathbb{KL}(p(x, z) \| p_\theta(x, z))$

Objective function:

$$\mathbb{KL}(p(x, z) \| p_\theta(x, z)) = \iint p(x)q_\phi(z | x) \log \frac{p(x)q_\phi(z | x)}{p_\theta(x | z)q(z)} dx dz$$

Alternating Optimization

1. Fix $q_\phi(z | x)$, optimize $p_\theta(x | z)$

We have

$$p_\theta(x | z) = \arg \max_{p_\theta(x|z)} \mathbb{E}_{x \sim p(x)} \left[\int q_\phi(z | x) \log p(x, z) dz \right] \quad (15)$$

2. ...

Alternating Algorithm

$$\mathbb{KL}(p(x, z) \| p_{\theta}(x, z)) = \iint p(x)q_{\phi}(z | x) \log \frac{p(x)q_{\phi}(z|x)}{p_{\theta}(x|z)q(z)} dx dz$$

Alternating Algorithm

$$\mathbb{KL}(p(x, z) \| p_\theta(x, z)) = \iint p(x)q_\phi(z | x) \log \frac{p(x)q_\phi(z | x)}{p_\theta(x|z)q(z)} dx dz$$

Alternating Optimization

1. ...
2. Fix $p_\theta(x | z)$, so $p(x, z)$ is fixed, then optimize $q_\phi(z | x)$
define $p_\theta(x | z)q(z) = p_\theta(z | x)p_\theta(x)$

$$p_\theta(x) = \int p_\theta(x | z)q(z)dz, \quad p_\theta(z | x) = \frac{p_\theta(x | z)q(z)}{p_\theta(x)} \quad (16)$$

Then we have

$$\begin{aligned} q_\phi(z | x) &= \arg \min_{q_\phi(z|x)} \mathbb{E}_{x \sim p(x)} \left[\int q_\phi(z | x) \log \frac{q_\phi(z | x)}{p_\theta(z | x)p_\theta(x)} dz \right] \\ &= \arg \min_{q_\phi(z|x)} \mathbb{E}_{x \sim p(x)} [\mathbb{KL}(q_\phi(z | x) \| p_\theta(z | x)) - \log p_\theta(x)] \end{aligned} \quad (17)$$

Alternating Algorithm

$q_\phi(z \mid x)$ have optimal solution:

Alternating Algorithm

$q_\phi(z | x)$ have optimal solution:

$$q_\phi(z | x) = p_\theta(z | x) = \frac{p_\theta(x | z)q(z)}{\int p_\theta(x | z)q(z)dz} \quad (18)$$

Alternating Algorithm

$q_\phi(z \mid x)$ have optimal solution:

$$q_\phi(z \mid x) = p_\theta(z \mid x) = \frac{p_\theta(x \mid z)q(z)}{\int p_\theta(x \mid z)q(z)dz} \quad (18)$$

$$\begin{aligned} p_\theta(x \mid z) &= \arg \max_{p_\theta(x \mid z)} \mathbb{E}_{x \sim p(x)} \left[\int q_\phi(z \mid x) \log p(x, z) dz \right] \\ &= \arg \max_{\theta} \mathbb{E}_{z \sim p_\theta(z \mid x)} [\log p(x, z)] \end{aligned}$$

Alternating Algorithm

$q_\phi(z | x)$ have optimal solution:

$$q_\phi(z | x) = p_\theta(z | x) = \frac{p_\theta(x | z)q(z)}{\int p_\theta(x | z)q(z)dz} \quad (18)$$

$$\begin{aligned} p_\theta(x | z) &= \arg \max_{p_\theta(x|z)} \mathbb{E}_{x \sim p(x)} \left[\int q_\phi(z | x) \log p(x, z) dz \right] \\ &= \arg \max_{\theta} \mathbb{E}_{z \sim p_\theta(z|x)} [\log p(x, z)] \end{aligned}$$

We get EM Algorithm by alternately computing these two formulae.

Diffusion Model: Autoregressive VAE

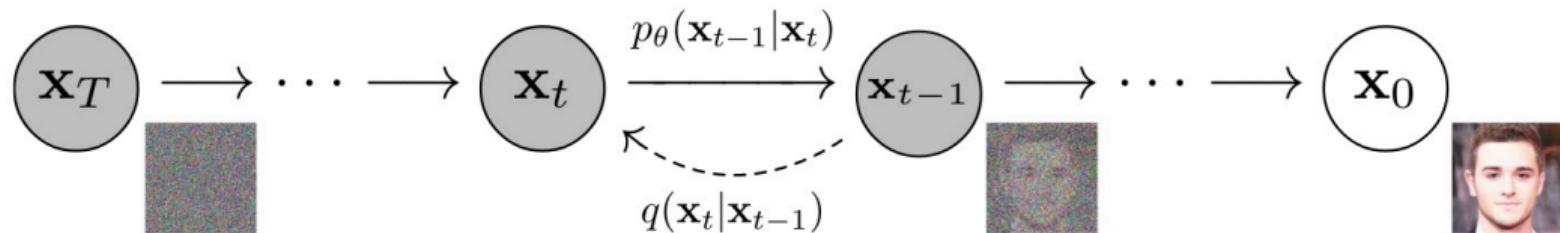


Figure 2: The directed graphical model considered in this work.

- Forward process: From right to left, model x_0, x_1, \dots, x_T as a Markov Chain:
$$x_t = \alpha_t x_{t-1} + \beta_t \quad \text{where} \quad \alpha_t^2 + \beta_t^2 = 1$$
- Reverse process: From left to right, learn $p_\theta(x_{t-1} | x_t)$ to fit $q(x_{t-1} | x_t)$.

Diffusion Model: Autoregressive VAE

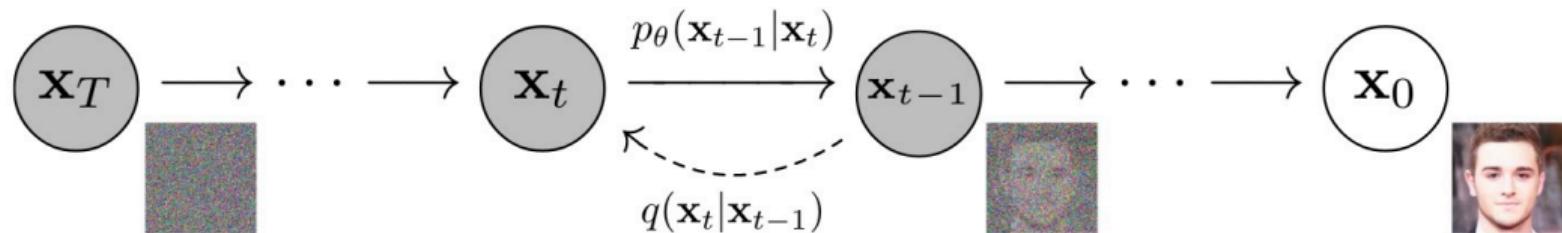


Figure 2: The directed graphical model considered in this work.

- Forward process: From right to left, model x_0, x_1, \dots, x_T as a Markov Chain:
$$x_t = \alpha_t x_{t-1} + \beta_t \quad \text{where} \quad \alpha_t^2 + \beta_t^2 = 1$$
- Reverse process: From left to right, learn $p_\theta(x_{t-1} | x_t)$ to fit $q(x_{t-1} | x_t)$.

$q(x_{t-1} | x_t, x_0)$ **can be represented by** $q(x_t | x_0)$ **and** $q(x_t | x_{t-1})$.

$q(x_{t-1} | x_t, x_0)$ **can be used to learn** $p_\theta(x_{t-1} | x_t)$

Diffusion Model: Autoregressive VAE

VAE

- encode: $x \rightarrow z$
- decode: $z \rightarrow x$

Diffusion Model: Autoregressive VAE

VAE

- encode: $x \rightarrow z$
- decode: $z \rightarrow x$

Diffusion Model

- encode: $x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{T-1} \rightarrow x_T = z$
- decode: $z = x_T \rightarrow x_{T-1} \rightarrow x_{T-2} \rightarrow \dots \rightarrow x_1 \rightarrow x_0 = x$

Diffusion Model: Autoregressive VAE

VAE

- encode: $x \rightarrow z$
- decode: $z \rightarrow x$

Diffusion Model

- encode: $x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{T-1} \rightarrow x_T = z$
- decode: $z = x_T \rightarrow x_{T-1} \rightarrow x_{T-2} \rightarrow \dots \rightarrow x_1 \rightarrow x_0 = x$

Drawback for VAE: encoder and decoder are both Gaussian distributions, the model is not complicated enough.

Diffusion Model: Autoregressive VAE

VAE

- encode: $x \rightarrow z$
- decode: $z \rightarrow x$

Diffusion Model

- encode: $x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{T-1} \rightarrow x_T = z$
- decode: $z = x_T \rightarrow x_{T-1} \rightarrow x_{T-2} \rightarrow \dots \rightarrow x_1 \rightarrow x_0 = x$

Drawback for VAE: encoder and decoder are both Gaussian distributions, the model is not complicated enough.

By modeling each step as a Gaussian distribution, the model becomes more complicated intuitively.

Model the Joint Distributions

For forward process, we define $q(x_t | x_{t-1})$

For reverse process, we define $p_\theta(x_{t-1} | x_t)$

Model the Joint Distributions

For forward process, we define $q(x_t | x_{t-1})$

For reverse process, we define $p_\theta(x_{t-1} | x_t)$

For Markov Chain assumption, we have

$$\begin{aligned} q(x_0, x_1, x_2, \dots, x_T) &= q(x_T | x_{T-1}) \cdots q(x_2 | x_1) q(x_1 | x_0) \tilde{q}(x_0) \\ p_\theta(x_0, x_1, x_2, \dots, x_T) &= p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-2} | x_{T-1}) p_\theta(x_{T-1} | x_T) p_\theta(x_T) \end{aligned} \tag{19}$$

Model the Joint Distributions

For forward process, we define $q(x_t | x_{t-1})$

For reverse process, we define $p_\theta(x_{t-1} | x_t)$

For Markov Chain assumption, we have

$$\begin{aligned} q(x_0, x_1, x_2, \dots, x_T) &= q(x_T | x_{T-1}) \cdots q(x_2 | x_1) q(x_1 | x_0) \tilde{q}(x_0) \\ p_\theta(x_0, x_1, x_2, \dots, x_T) &= p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-2} | x_{T-1}) p_\theta(x_{T-1} | x_T) p_\theta(x_T) \end{aligned} \tag{19}$$

We use the KL divergence between two joint distributions as our objective.

$$\text{KL}(q \| p_\theta) = \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log \frac{q(x_T | x_{T-1}) \cdots p(x_1 | x_0) \tilde{q}(x_0)}{p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T)} dx_0 \cdots dx_T \tag{20}$$

Derive the KL divergence

$$\begin{aligned} & \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log \frac{q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0)}{p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T)} dx_0 dx_1 \cdots dx_T \\ &= - \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T) dx_0 dx_1 \cdots dx_T \\ &= - \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \left[\log p_\theta(x_T) + \sum_{t=1}^T \log p_\theta(x_{t-1} | x_t) \right] dx_0 dx_1 \cdots dx_T \end{aligned} \tag{21}$$

Derive the KL divergence

$$\begin{aligned} & \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log \frac{q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0)}{p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T)} dx_0 dx_1 \cdots dx_T \\ &= - \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T) dx_0 dx_1 \cdots dx_T \\ &= - \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \left[\log p_\theta(x_T) + \sum_{t=1}^T \log p_\theta(x_{t-1} | x_t) \right] dx_0 dx_1 \cdots dx_T \end{aligned} \tag{21}$$

Typically, $p_\theta(x_T)$ is defined as standard Gaussian distribution.

For each $p_\theta(x_{t-1} | x_t)$, we have

$$-\int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_T \tag{22}$$

Derive the KL divergence

$$\begin{aligned} & \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log \frac{q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0)}{p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T)} dx_0 dx_1 \cdots dx_T \\ &= - \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_0 | x_1) \cdots p_\theta(x_{T-1} | x_T) p_\theta(x_T) dx_0 dx_1 \cdots dx_T \\ &= - \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \left[\log p_\theta(x_T) + \sum_{t=1}^T \log p_\theta(x_{t-1} | x_t) \right] dx_0 dx_1 \cdots dx_T \end{aligned} \tag{21}$$

Typically, $p_\theta(x_T)$ is defined as standard Gaussian distribution.

For each $p_\theta(x_{t-1} | x_t)$, we have

$$- \int q(x_T | x_{T-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_T \tag{22}$$

$p_\theta(x_{t-1}, x_t)$ is only related to $x_0 \cdots x_t$, so

$$= - \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_t \tag{23}$$

Derive the KL divergence

$$= - \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_t$$

Derive the KL divergence

$$= - \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_t$$

According to the following proposition, we have

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t \quad (24)$$

Derive the KL divergence

$$= - \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_t$$

According to the following proposition, we have

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t \quad (24)$$

Proposition

Given by [Ho, Neurips 2020], we have

$$q(x_t | x_0) = \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) dx_1 \cdots dx_{t-1} = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 I) \quad (25)$$

where $x_t = \alpha_t x_{t-1} + \beta_t \varepsilon_t$, $\varepsilon_t \sim \mathcal{N}(\mathbf{0}, I)$, $\bar{\alpha}_t = \alpha_1 \cdots \alpha_t$, $\bar{\beta}_t = \sqrt{1 - (\alpha_t \cdots \alpha_1)^2}$

Derive the KL divergence

$$= - \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_1 \cdots dx_t$$

According to the following proposition, we have

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t \quad (24)$$

Proposition

Given by [Ho, Neurips 2020], we have

$$q(x_t | x_0) = \int q(x_t | x_{t-1}) \cdots q(x_1 | x_0) dx_1 \cdots dx_{t-1} = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 I) \quad (25)$$

where $x_t = \alpha_t x_{t-1} + \beta_t \varepsilon_t$, $\varepsilon_t \sim \mathcal{N}(\mathbf{0}, I)$, $\bar{\alpha}_t = \alpha_1 \cdots \alpha_t$, $\bar{\beta}_t = \sqrt{1 - (\alpha_t \cdots \alpha_1)^2}$

So $p_\theta(x_{t-1} | x_t)$ is irrelevant to x_1, \dots, x_{t-2} and we have (24).

Objective Function

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t$$

1. Ignore those constants irrelevant to optimization objective, $-\log p_\theta(x_{t-1} | x_t)$ is $\frac{1}{2\sigma_t^2} \|x_{t-1} - \mu(x_t)\|^2$

Objective Function

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t$$

1. Ignore those constants irrelevant to optimization objective, $-\log p_\theta(x_{t-1} | x_t)$ is $\frac{1}{2\sigma_t^2} \|x_{t-1} - \mu(x_t)\|^2$
2. With $x_{t-1} = \frac{1}{\alpha_t} (x_t - \beta_t \varepsilon_t)$, we can parameterize $\mu(x_t)$ as
$$\mu(x_t) = \frac{1}{\alpha_t} (x_t - \beta_t \epsilon_\theta(x_t, t))$$

Objective Function

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t$$

1. Ignore those constants irrelevant to optimization objective, $-\log p_\theta(x_{t-1} | x_t)$ is $\frac{1}{2\sigma_t^2} \|x_{t-1} - \mu(x_t)\|^2$
2. With $x_{t-1} = \frac{1}{\alpha_t} (x_t - \beta_t \varepsilon_t)$, we can parameterize $\mu(x_t)$ as
$$\mu(x_t) = \frac{1}{\alpha_t} (x_t - \beta_t \epsilon_\theta(x_t, t))$$
3. According to our proposition, $q(x_{t-1} | x_0)$ is $x_{t-1} = \bar{\alpha}_{t-1} x_0 + \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1}$

Objective Function

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t$$

1. Ignore those constants irrelevant to optimization objective, $-\log p_\theta(x_{t-1} | x_t)$ is $\frac{1}{2\sigma_t^2} \|x_{t-1} - \mu(x_t)\|^2$
2. With $x_{t-1} = \frac{1}{\alpha_t} (x_t - \beta_t \varepsilon_t)$, we can parameterize $\mu(x_t)$ as
$$\mu(x_t) = \frac{1}{\alpha_t} (x_t - \beta_t \epsilon_\theta(x_t, t))$$
3. According to our proposition, $q(x_{t-1} | x_0)$ is $x_{t-1} = \bar{\alpha}_{t-1} x_0 + \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1}$
4. $q(x_t | x_{t-1})$ is $x_t = \alpha_t x_{t-1} + \beta_t \varepsilon_t$, where $\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{0}, I)$

Objective Function

$$= - \int q(x_t | x_{t-1}) q(x_{t-1} | x_0) \tilde{q}(x_0) \log p_\theta(x_{t-1} | x_t) dx_0 dx_{t-1} dx_t$$

1. Ignore those constants irrelevant to optimization objective, $-\log p_\theta(x_{t-1} | x_t)$ is $\frac{1}{2\sigma_t^2} \|x_{t-1} - \mu(x_t)\|^2$
2. With $x_{t-1} = \frac{1}{\alpha_t} (x_t - \beta_t \varepsilon_t)$, we can parameterize $\mu(x_t)$ as
$$\mu(x_t) = \frac{1}{\alpha_t} (x_t - \beta_t \epsilon_\theta(x_t, t))$$
3. According to our proposition, $q(x_{t-1} | x_0)$ is $x_{t-1} = \bar{\alpha}_{t-1} x_0 + \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1}$
4. $q(x_t | x_{t-1})$ is $x_t = \alpha_t x_{t-1} + \beta_t \varepsilon_t$, where $\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{0}, I)$

$$= \frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{0}, I), x_0 \sim \tilde{q}(x_0)} \left[\|\varepsilon_t - \epsilon_\theta(\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t, t)\|^2 \right] \quad (26)$$

Objective Function

$$\frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(0, I), x_0 \sim \tilde{q}(x_0)} \left[\|\varepsilon_t - \epsilon_\theta (\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t, t)\|^2 \right]$$

Objective Function

$$\frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{q}(x_0)} \left[\|\varepsilon_t - \epsilon_\theta (\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t, t)\|^2 \right]$$

For this equation, we have $\bar{\varepsilon}_{t-1}, \varepsilon_t$ to randomly generate, we use a trick [Ho, Neurips 2020] to replace them with one random variable ε .

Proposition

$\alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t$ is equivalent to $\bar{\beta}_t \varepsilon \mid \varepsilon \sim \mathcal{N}(\mathbf{o}, I)$

$\beta_t \bar{\varepsilon}_{t-1} - \alpha_t \bar{\beta}_{t-1} \varepsilon_t$ is equivalent to $\bar{\beta}_t \omega \mid \omega \sim \mathcal{N}(\mathbf{o}, I)$

It can be verified that $\mathbb{E} [\varepsilon \omega^\top] = \mathbf{o}$

Objective Function

$$\frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{q}(x_0)} \left[\|\varepsilon_t - \epsilon_\theta (\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t, t)\|^2 \right]$$

For this equation, we have $\bar{\varepsilon}_{t-1}, \varepsilon_t$ to randomly generate, we use a trick [Ho, Neurips 2020] to replace them with one random variable ε .

Proposition

$\alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t$ is equivalent to $\bar{\beta}_t \varepsilon$ | $\varepsilon \sim \mathcal{N}(\mathbf{o}, I)$

$\beta_t \bar{\varepsilon}_{t-1} - \alpha_t \bar{\beta}_{t-1} \varepsilon_t$ is equivalent to $\bar{\beta}_t \omega$ | $\omega \sim \mathcal{N}(\mathbf{o}, I)$

It can be verified that $\mathbb{E} [\varepsilon \omega^\top] = \mathbf{o}$

We use ε and ω to represent ε_t

$$\varepsilon_t = \frac{(\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega) \bar{\beta}_t}{\beta_t^2 + \alpha_t^2 \bar{\beta}_{t-1}^2} = \frac{\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t} \quad (27)$$

Objective Function

$$\varepsilon_t = \frac{(\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega) \bar{\beta}_t}{\beta_t^2 + \alpha_t^2 \bar{\beta}_{t-1}^2} = \frac{\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t}$$

Objective Function

$$\varepsilon_t = \frac{(\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega) \bar{\beta}_t}{\beta_t^2 + \alpha_t^2 \bar{\beta}_{t-1}^2} = \frac{\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t}$$

Plug this equation into our objective.

$$\begin{aligned} & \frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{q}(x_0)} \left[\left\| \varepsilon_t - \epsilon_\theta (\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t, t) \right\|^2 \right] \\ &= \frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\omega, \varepsilon \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{q}(x_0)} \left[\left\| \frac{\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t} - \epsilon_\theta (\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t) \right\|^2 \right] \end{aligned} \tag{28}$$

Objective Function

$$\varepsilon_t = \frac{(\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega) \bar{\beta}_t}{\beta_t^2 + \alpha_t^2 \bar{\beta}_{t-1}^2} = \frac{\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t}$$

Plug this equation into our objective.

$$\begin{aligned} & \frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\varepsilon}_{t-1}, \varepsilon_t \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{q}(x_0)} \left[\left\| \varepsilon_t - \epsilon_\theta (\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \bar{\varepsilon}_{t-1} + \beta_t \varepsilon_t, t) \right\|^2 \right] \\ &= \frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\omega, \varepsilon \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{q}(x_0)} \left[\left\| \frac{\beta_t \varepsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t} - \epsilon_\theta (\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t) \right\|^2 \right] \end{aligned} \tag{28}$$

Notice that ω can be separated in this equation, we have

$$\frac{\beta_t^4}{\bar{\beta}_t^2 \alpha_t^2 \sigma_t^2} \mathbb{E}_{\varepsilon \sim \mathcal{N}(\mathbf{o}, I), x_0 \sim \tilde{p}(x_0)} \left[\left\| \varepsilon - \frac{\bar{\beta}_t}{\beta_t} \epsilon_\theta (\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t) \right\|^2 \right] \tag{29}$$

Score-based Methods for Diffusion Model

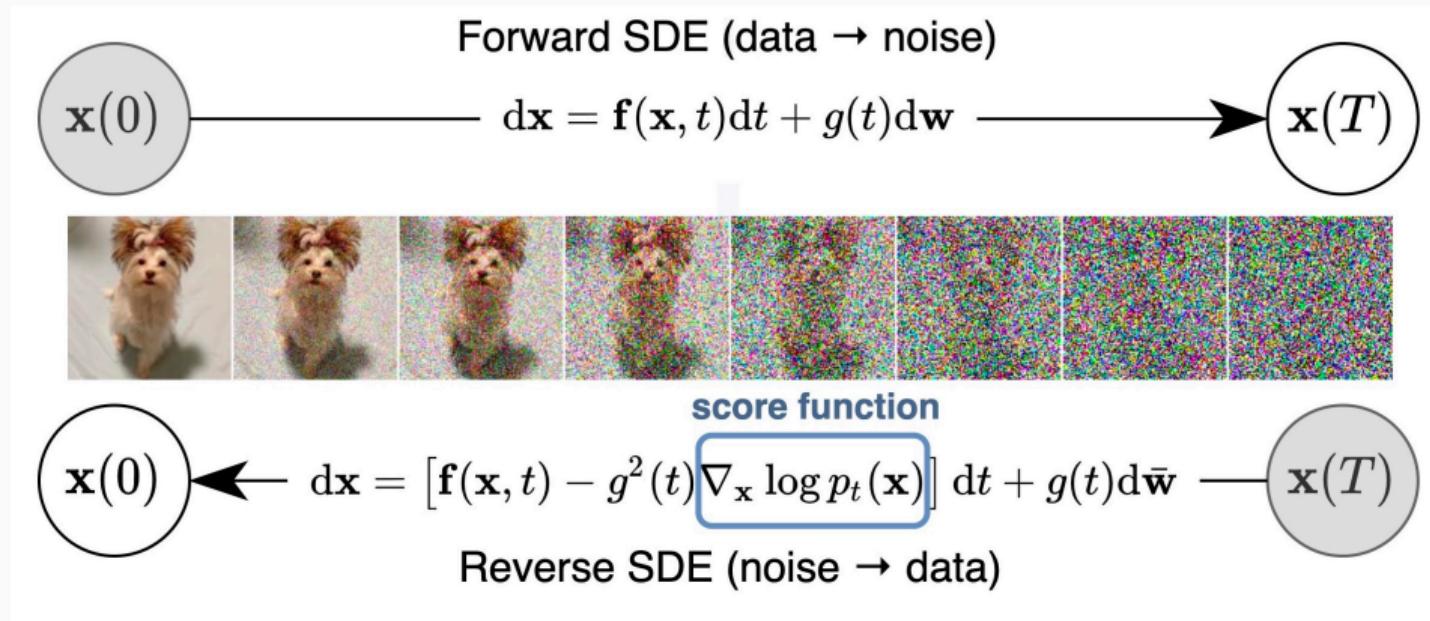


Figure 3: Solving a reverse SDE yields a score-based generative model.

note: $p_t(x) = p(x_t)$, where $x_t = \text{func}(x, t)$

Connection to DDPM

DDPM:

$$x_i = \sqrt{1 - \beta_i} x_{i-1} + \sqrt{\beta_i} z_{i-1}, \quad i = 1, \dots, N \quad (30)$$

Connection to DDPM

DDPM:

$$x_i = \sqrt{1 - \beta_i} x_{i-1} + \sqrt{\beta_i} z_{i-1}, \quad i = 1, \dots, N \quad (30)$$

Corresponding Stochastic Differential Equation:

$$dx = -\frac{1}{2}\beta(t)xdt + \sqrt{\beta(t)}dw \quad (31)$$

Connection to DDPM

DDPM:

$$x_i = \sqrt{1 - \beta_i} x_{i-1} + \sqrt{\beta_i} z_{i-1}, \quad i = 1, \dots, N \quad (30)$$

Corresponding Stochastic Differential Equation:

$$dx = -\frac{1}{2}\beta(t)xdt + \sqrt{\beta(t)}dw \quad (31)$$

DDPM is one of the discretization version of SDE, with specific coefficients.

Solve the Reverse SDE

Reverse SDE:

$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)] dt + g(t) d\bar{w} \quad (32)$$

Solve the Reverse SDE

Reverse SDE:

$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)] dt + g(t) d\bar{w} \quad (32)$$

Unknown distributions:

- the terminal distribution $p_T(x) \approx \pi(x)$
- the score function $\nabla_x \log p_t(x)$

In order to estimate $\nabla_x \log p_t(x)$, we train a Time-Dependent Score-Based Model $s_\theta(x, t)$, such that $s_\theta(x, t) \approx \nabla_x \log p_t(x)$

$$\mathbb{E}_{t \in \mathcal{U}(0, T)} \mathbb{E}_{p_t(x)} [\lambda(t) \|\nabla_x \log p_t(x) - s_\theta(x, t)\|_2^2] \quad (33)$$

Solve the Reverse SDE

Reverse SDE:

$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)] dt + g(t) d\bar{w} \quad (32)$$

Unknown distributions:

- the terminal distribution $p_T(x) \approx \pi(x)$
- the score function $\nabla_x \log p_t(x)$

In order to estimate $\nabla_x \log p_t(x)$, we train a Time-Dependent Score-Based Model $s_\theta(x, t)$, such that $s_\theta(x, t) \approx \nabla_x \log p_t(x)$

$$\mathbb{E}_{t \in \mathcal{U}(0, T)} \mathbb{E}_{p_t(x)} [\lambda(t) \|\nabla_x \log p_t(x) - s_\theta(x, t)\|_2^2] \quad (33)$$

There exists a family of methods called score matching [Hyvarinen JMLR 2005; Song, UAI 2020] that minimize the divergence without knowledge of the ground-truth score function.

Connection to Likelihood Estimation

Proposition

Given by [Song, ICLR 2021]:

When $\lambda(t) = g(t)^2$, we have

$$\mathbb{KL}(p_0(x) \| p_\theta(x)) \leq \frac{T}{2} \mathbb{E}_{t \in \mathcal{U}(0, T)} \mathbb{E}_{p_t(x)} [\lambda(t) \|\nabla_x \log p_t(x) - s_\theta(x, t)\|_2^2] + \mathbb{KL}(p_T \| \pi) \quad (34)$$

According to this, we can train score-based generative models to achieve very high likelihoods.

Take Home Message

- **Variational inference can be used to derive many likelihood-based generative models including VAE, DDPM, etc.**
- **ELBO is a lower bound for $\log p_\theta(x)$**

$$\Updownarrow$$

$\mathbb{KL}(p(x, z) \| p_\theta(x, z))$ **is an upper bound for $\mathbb{KL}(p(x) \| p_\theta(x))$**

- **DDPM is an autoregressive version of VAE and one of the discretization version of SDE, with specific coefficients.**

Thanks!