

扩散生成模型 Diffusion Model理论基础

References:

《Understanding Diffusion Models: A Unified Perspective》

<https://arxiv.org/abs/2208.11970>

《Denoising Diffusion Probabilistic Models》

<https://arxiv.org/abs/2006.11239>

目录

- VAE数学原理及模型结构
- MHVAE理论推导
- VDM理论推导(DDPM)

VAE数学原理及模型结构

Variational AutoEncoder

VAE(Variational Autoencoders)简介

Likelihood-based生成模型：给定一个数据集 x_D ，训练使得模型最大化likelihood $p_\phi(x_D)$

$$x_D \xrightarrow[\phi]{\text{model}} p_\phi(x) \xrightarrow{\text{sample}} x'$$

VAE：生成需要采样，借助一个变量 z 和自定义分布 $p(z)$ ，一般选择多维标准高斯分布 $N \sim (z; \mathbf{0}, \mathbf{I})$

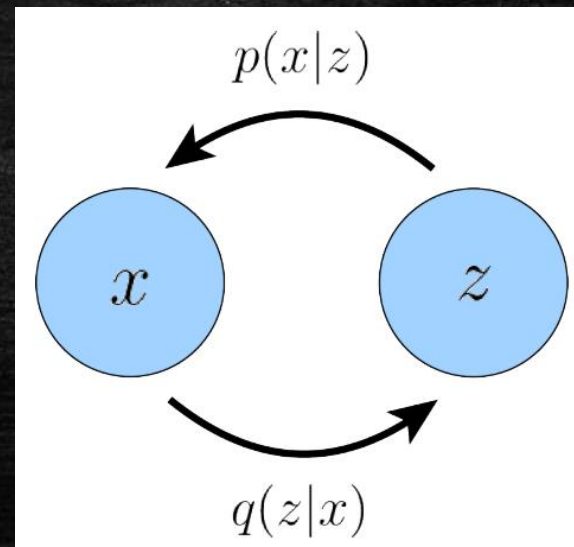
$q(z|x)$ ：Encoder

$p(x|z)$ ：Decoder

将 z 和 x 建立起了联系

对高斯分布的 z 采样，就能通过decoder得到一个新的生成数据 x'

z ：latent variable



VAE数学原理

VAE模型优化目标: $p_\phi(x)$ 趋近于真实数据分布 $p_\phi(x) \rightarrow p(x)$

$\log P(x)$ 有下界:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \quad (1)$$

ELBO:

Evidence

Lower Bound

ELBO公式证明方法一 (琴生不等式):

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

边缘化

$$= \log \int \frac{p(\mathbf{x}, \mathbf{z}) q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$\times 1$

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

期望的定义

$$\mathbb{E}_{P(B)} \left[\frac{p(A)}{p(B)} \right] = \int \frac{p(A)}{p(B)} \cdot p(B) dB$$

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

琴生不等式

https://upload.wikimedia.org/wikipedia/commons/transcoded/5/52/Convex_01.ogv/Convex_01.ogv.360p.webm

VAE数学原理

ELBO是 $\log P(x)$ 下界，但他们之间具体有什么关系呢？
什么时候能取到等号？
为什么说： **优化VAE \Leftrightarrow 最大化ELBO**

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

ELBO公式证明方法二：

$$\begin{aligned} \log p(\mathbf{x}) &= \log p(\mathbf{x}) \int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) (\log p(\mathbf{x})) d\mathbf{z} \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) \\ &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \end{aligned}$$

$\times 1$, 任一分布概率密度函数的积分等于1 $\leftarrow \int p(A) dA = 1$

期望的定义 $\leftarrow E_{P(B)}[p(A)] = \int p(A) \cdot p(B) dB$

链式法则 $\leftarrow p(x, z) = p(x) \cdot p(z|x)$

$\times 1$

和的期望=期望的和 $\leftarrow E[\log(AB)] = E[\log(A)] + E[\log(B)]$

KL散度定义 $\leftarrow D_{\text{KL}}[p(A) \parallel p(B)] = E_{P(A)} \left[\log \frac{p(A)}{p(B)} \right]$

KL散度 ≥ 0

VAE数学原理

$$\log p(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$

ELBO是 $\log P(x)$ 下界，但他们之间具体有什么关系呢？

答：差1个KL散度，模型encoder拟合的 $q_{\phi}(z|x)$ 和真实的 $p(z|x)$ 之间的KL散度

什么时候能取到等号？

答：KL散度=0，模型encoder完美拟合真实的 $p(z|x)$ ，及 $q_{\phi}(z|x) = p(z|x)$

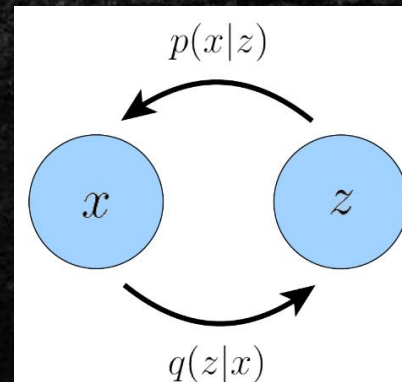
为什么说：优化VAE \Leftrightarrow 最大化ELBO？

答：VAE的模型框架就是拟合数据变量 x 和隐变量 z 之间的联系。

也就是通过decoder和encoder分别拟合 $p(x|z)$ 和 $q(z|x)$ ，间接拟合 $p_{\phi}(x) \rightarrow p(x)$

KL散度项反映的就是encoder对于 $p(z|x)$ 的拟合程度，我们要最小化它。但是直接优化这个KL散度不可行，因为我们不知道 $p(z|x)$ 这个groundtruth。

$\log p(x)$ 是真实数据分布的概率，给定一个数据集，该值就唯一确定，与模型无关。
所以：最小化KL散度 \Leftrightarrow 最大化ELBO。



VAE数学原理

我们继续拆解ELBO，看看其有什么具体含义？

$$\begin{aligned}\mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(\mathbf{x}, z)}{q_\phi(z|x)} \right] &= \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(\mathbf{x}|z)p(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(\mathbf{x}|z)] + \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(z)}{q_\phi(z|x)} \right] \\ &= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(\mathbf{x}|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}\end{aligned}$$

链式法则

$$\leftarrow p(\mathbf{x}, z) = p(\mathbf{x}) \cdot p(z|x)$$

和的期望=期望的和

KL散度定义

$$\leftarrow D_{\text{KL}}[p(A) \parallel p(B)] = E_{P(A)} \left[\log \frac{p(A)}{p(B)} \right]$$

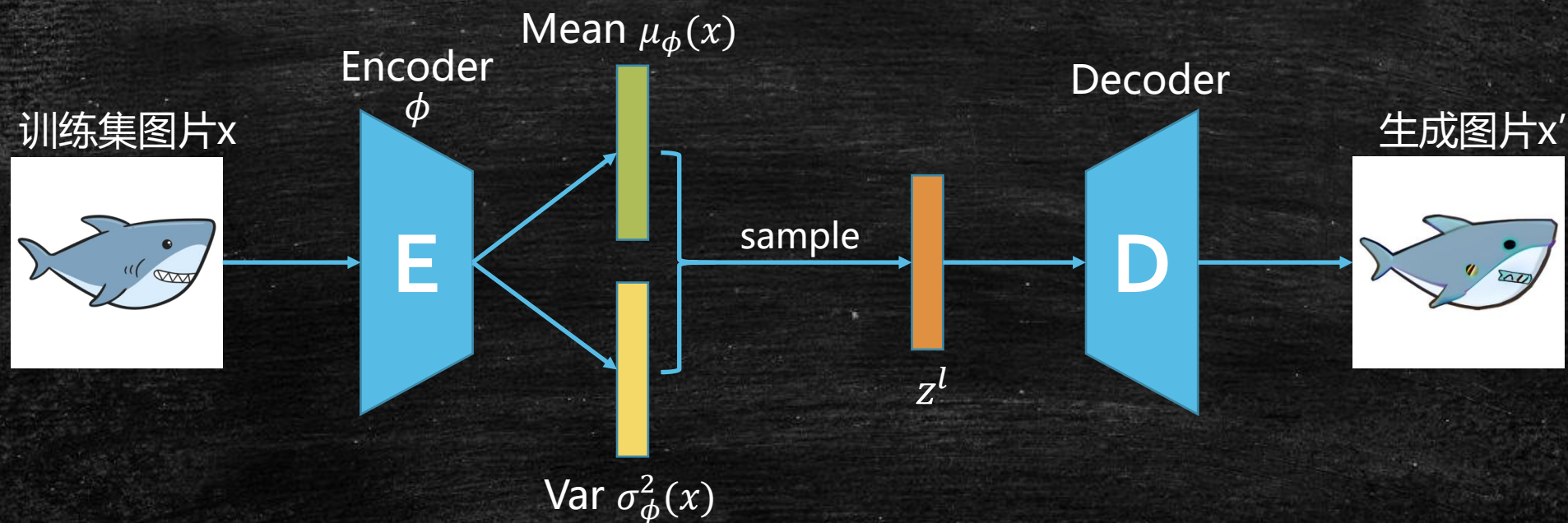
最大化**Reconstruction term**：让decoder能最大可能从隐变量 z 生成原始的真实数据 x

最小化**Prior matching term**：让encoder将真实数据 x 映射到隐变量 z 后， z 尽量满足我们指定的分布，一般为多维标准高斯分布 $N \sim (z; \mathbf{0}, I)$

TIPs：如果把**Prior matching term**这项损失去掉，就是AE模型，AE的 z 分布未知，没法有效采样，所以也不能作为生成模型

VAE模型结构

VAE模型结构具体是怎么实现的？与ELBO的两部分对应关系是什么？



最大化**Reconstruction term**：生成图片 x' 和训练集图片 x 尽可能趋近

最小化**Prior matching term**： $q_\phi(z|x) = N(z; \mu_\phi(x), \sigma_\phi^2(x)I) \rightarrow p(z) = N(z; 0, I)$

KL散度趋近于0， mean向量和var向量分别趋近于全0向量和全1向量

VAE模型结构

重参数化技巧是什么？有什么好处？

Mean $\mu_\phi(x)$



sample

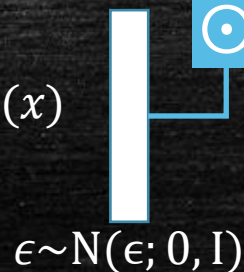


重参数化

Mean $\mu_\phi(x)$



Var $\sigma_\phi^2(x)$



+

\odot



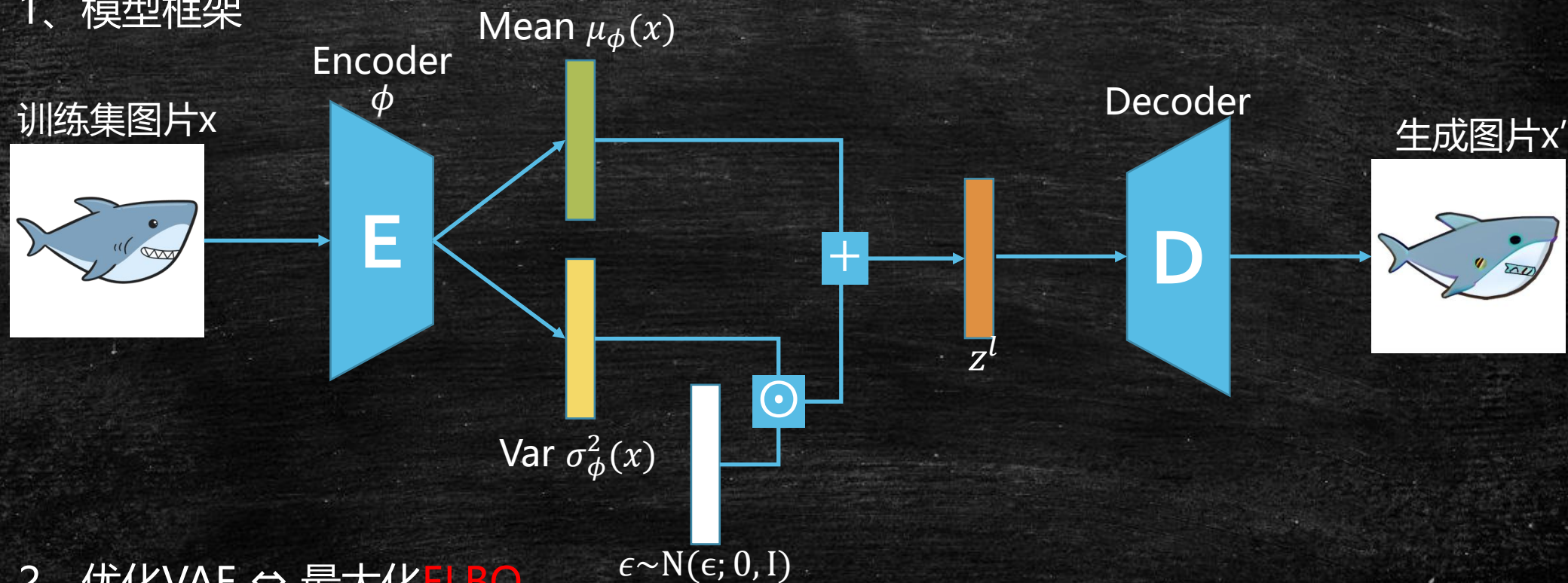
$z^l = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon$
 \odot 代表元素相乘，这种计算下不改变 z 的分布

这样参数 ϕ 被剥离出了随机采样过程，含参的随机采样变成了不含参采样，这就是重参数化过程，这样参数 ϕ 可导

$z \sim N(z; \mu_\phi(x), \sigma_\phi^2(x)I)$ 分布采样是随机过程，随机过程中包含要优化的参数 ϕ ，但是随机采样过程对参数 ϕ 不可导

VAE模型小结

1、模型框架



2、优化VAE \Leftrightarrow 最大化ELBO

3、ELBO包含两部分，一部分Reconstruction term反映decoder从隐变量重建图片的能力，一部分Prior matching term反映encoder将图片映射到指定隐变量分布的能力

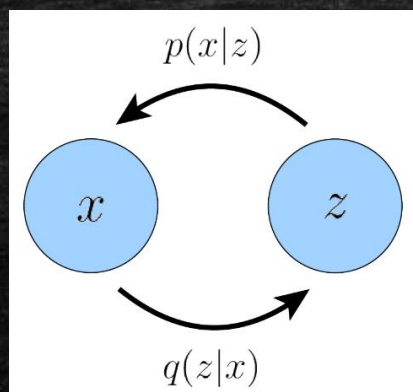
4、重参数化技巧将模型参数剥离出随机采样过程，使其可导

5、生成时，只需要从高斯分布随机采样一个因变量 z ，经过decoder，就能生成图片啦

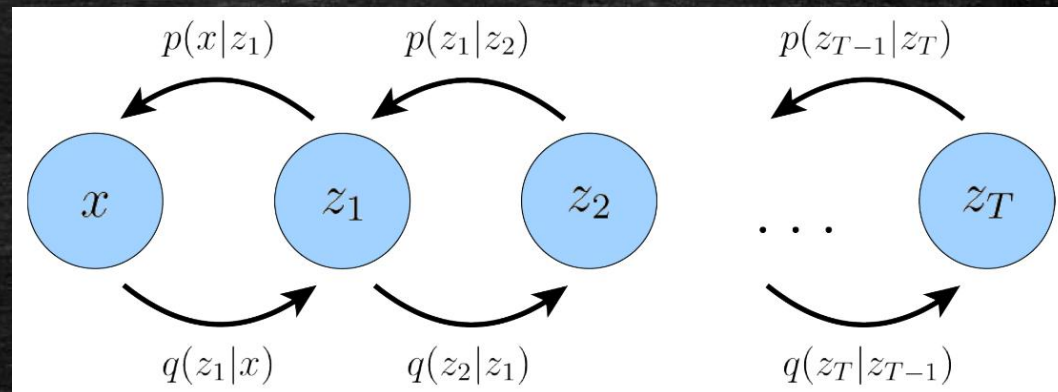
MHVAE理论推导

Markovian Hierarchical Variational AutoEncoder

从VAE到MHVAE



级联
Hierarchical
马尔科夫链
Markovian



$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{z}) q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$

边缘化

$\times 1$

期望的定义

琴生不等式

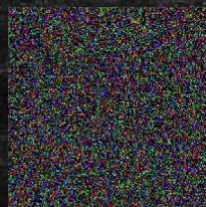
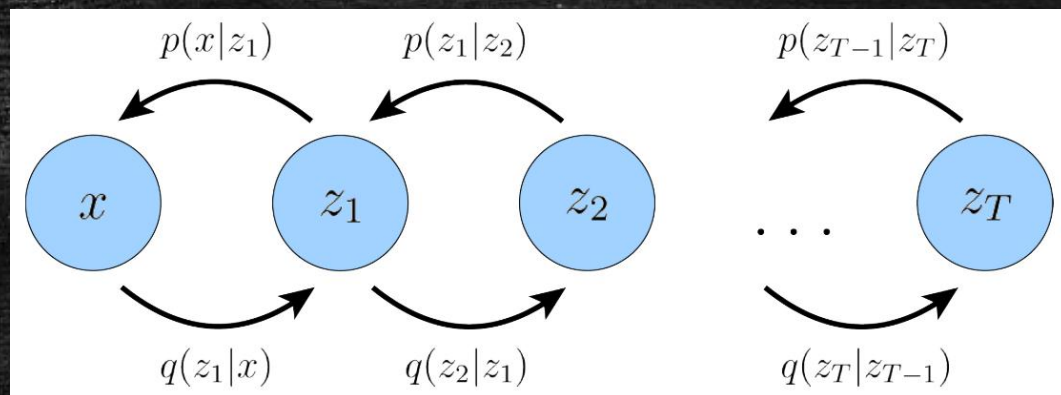
$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{z}_{1:T}) q_\phi(\mathbf{z}_{1:T}|\mathbf{x})}{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} d\mathbf{z}_{1:T} \\ &= \log \mathbb{E}_{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \right] \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \right]\end{aligned}$$

VDM理论推导

Variational Diffusion Models

从MHVAE到VDM

MHVAE



三个限制:

- ✓ 数据 x 和所有隐变量 z_t 的维度相同
- ✓ 所有的encoder $q(z_t|z_{t-1})$ 都不需要学习, 而是预定义好的高斯分布模型, 就是 z_t 状态为以 z_{t-1} 为均值的高斯分布
- ✓ 最终 T 状态的分布 z_T 为标准高斯分布



VDM

从MHVAE到VDM

VDM三个限制:

- ① 数据 x 和所有隐变量 z_t 的维度相同
- ② 所有的encoder $q(z_t|z_{t-1})$ 都不需要学习, 而是预定义好的高斯分布模型, 就是 z_t 状态为以 z_{t-1} 为均值的高斯分布
- ③ 最终T状态的分布 z_T 为标准高斯分布

符号表示替换:

$$x \rightarrow x_0$$

$$z_t \rightarrow x_t$$



x_0



x_t



x_T

限制①②:

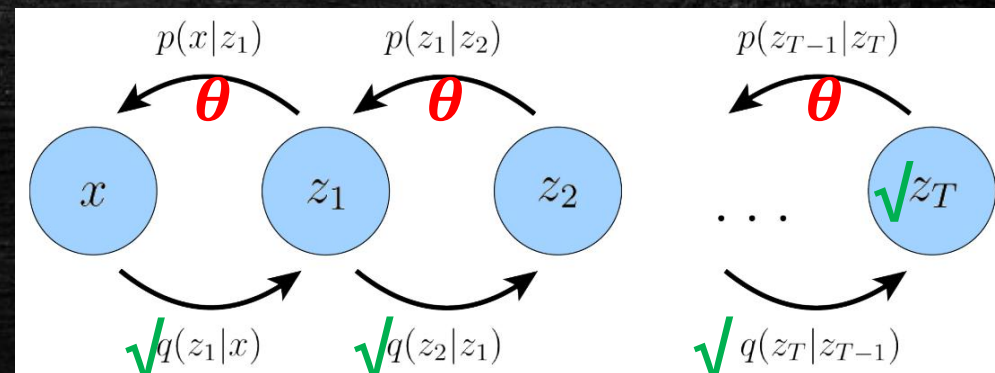
人为设定 $q(x_t|x_{t-1})$ 满足如下高斯分布

$$q(x_t|x_{t-1}) \sim N(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

α 是超参, 一般人为指定, 例如SD中的 noise schedule, 也可以通过模型学习

限制③:

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



VDM数学推导

VDM的ELBO证明
和MHVAE的证明只是把符号表示替换了一下

$$\begin{aligned} \mathbf{x} &\rightarrow \mathbf{x}_0 \\ \mathbf{z}_t &\rightarrow \mathbf{x}_t \end{aligned}$$

$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T}$	→	边缘化	→	$\log p(\mathbf{x}) = \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$
$= \log \int \frac{p(\mathbf{x}, \mathbf{z}_{1:T}) q_\phi(\mathbf{z}_{1:T} \mathbf{x})}{q_\phi(\mathbf{z}_{1:T} \mathbf{x})} d\mathbf{z}_{1:T}$	→	×1	→	$= \log \int \frac{p(\mathbf{x}_{0:T}) q(\mathbf{x}_{1:T} \mathbf{x}_0)}{q(\mathbf{x}_{1:T} \mathbf{x}_0)} d\mathbf{x}_{1:T}$
$= \log \mathbb{E}_{q_\phi(\mathbf{z}_{1:T} \mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T} \mathbf{x})} \right]$	→	期望的定义	→	$= \log \mathbb{E}_{q(\mathbf{x}_{1:T} \mathbf{x}_0)} \left[\frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mathbf{x}_0)} \right]$
$\geq \mathbb{E}_{q_\phi(\mathbf{z}_{1:T} \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T} \mathbf{x})} \right]$	→	琴生不等式	→	$\geq \mathbb{E}_{q(\mathbf{x}_{1:T} \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mathbf{x}_0)} \right]$

VDM数学推导

VDM的ELBO拆解，看其具体的含义(页1/2)

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]$$

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T)p(\mathbf{x}_{0:T-1}|\mathbf{x}_T)$$

$$= p(\mathbf{x}_T)p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_{0:T-2}|\mathbf{x}_{T-1}, \mathbf{x}_T)$$

$$= p(\mathbf{x}_T)p(\mathbf{x}_{T-1}|\mathbf{x}_T)p(\mathbf{x}_{0:T-2}|\mathbf{x}_{T-1})$$

\vdots

$$= p(\mathbf{x}_T)p(\mathbf{x}_{T-1}|\mathbf{x}_T) \dots p(\mathbf{x}_0|\mathbf{x}_1)$$

$$= p(\mathbf{x}_T) \prod_{t=1}^T p(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

链式
法则

马尔科夫性质

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = q(\mathbf{x}_{2:T}|\mathbf{x}_1)q(\mathbf{x}_1|\mathbf{x}_0)$$

$$= q(\mathbf{x}_{3:T}|\mathbf{x}_2, \mathbf{x}_1)q(\mathbf{x}_2|\mathbf{x}_1)q(\mathbf{x}_1|\mathbf{x}_0)$$

$$= q(\mathbf{x}_{3:T}|\mathbf{x}_2)q(\mathbf{x}_2|\mathbf{x}_1)q(\mathbf{x}_1|\mathbf{x}_0)$$

\vdots

$$= q(\mathbf{x}_T|\mathbf{x}_{T-1}) \dots q(\mathbf{x}_2|\mathbf{x}_1)q(\mathbf{x}_1|\mathbf{x}_0)$$

$$= \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

链式
法则

马尔科夫性质

VDM数学推导

VDM的ELBO拆解，看其具体的含义(页1/2)

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad \rightarrow \text{条件概率+马尔科夫性质} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad \rightarrow \text{分子分母连乘里各拆出一项} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=1}^{T-1} p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad \rightarrow \text{连乘部分的符号都减1} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \prod_{t=1}^{T-1} \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad \rightarrow \text{和的期望=期望的和} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^{T-1} \log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]\end{aligned}$$

VDM数学推导

VDM的ELBO拆解，看其具体的含义(页2/2)

$$\begin{aligned}
 \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\sum_{t=1}^{T-1} \log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \rightarrow \text{和的期望=期望的和} \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{T-1}, \mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \rightarrow \text{删掉无关的变量} \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))]}_{\text{prior matching term}} \\
 &\quad - \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}
 \end{aligned}$$

见下两页

VDM数学推导

VDM的ELBO拆解, Prior matching term最后一步推导:

$$\begin{aligned} & \mathbb{E}_{q(x_{T-1}, x_T | x_0)} \left[\log \frac{p(x_T)}{q(x_T | x_{T-1})} \right] \quad \boxed{\text{Prior matching term}} \\ &= \iint \left[\log \frac{p(x_T)}{q(x_T | x_{T-1})} \right] q(x_{T-1}, x_T | x_0) dx_{T-1} dx_T \quad \rightarrow \quad \boxed{\text{多元函数的期望定义}} \\ &= \iint \left[\log \frac{p(x_T)}{q(x_T | x_{T-1})} \right] q(x_{T-1} | x_0) dx_{T-1} q(x_T | x_{T-1}, x_0) dx_T \quad \rightarrow \quad \boxed{\text{链式法则}} \\ &= \int \left[\mathbb{E}_{q(x_T | x_{T-1})} \log \frac{p(x_T)}{q(x_T | x_{T-1})} \right] q(x_{T-1} | x_0) dx_{T-1} \quad \rightarrow \quad \boxed{\text{期望的定义+马尔科夫性质}} \\ &= - \int [D_{KL}(q(x_T | x_{T-1}) || p(x_T))] q(x_{T-1} | x_0) dx_{T-1} \quad \rightarrow \quad \boxed{\text{KL散度的定义}} \\ &= - \mathbb{E}_{q(x_{T-1} | x_0)} [D_{KL}(q(x_T | x_{T-1}) || p(x_T))] \quad \rightarrow \quad \boxed{\text{期望的定义}} \end{aligned}$$

VDM数学推导

VDM的ELBO拆解, Consistency term最后一步推导:

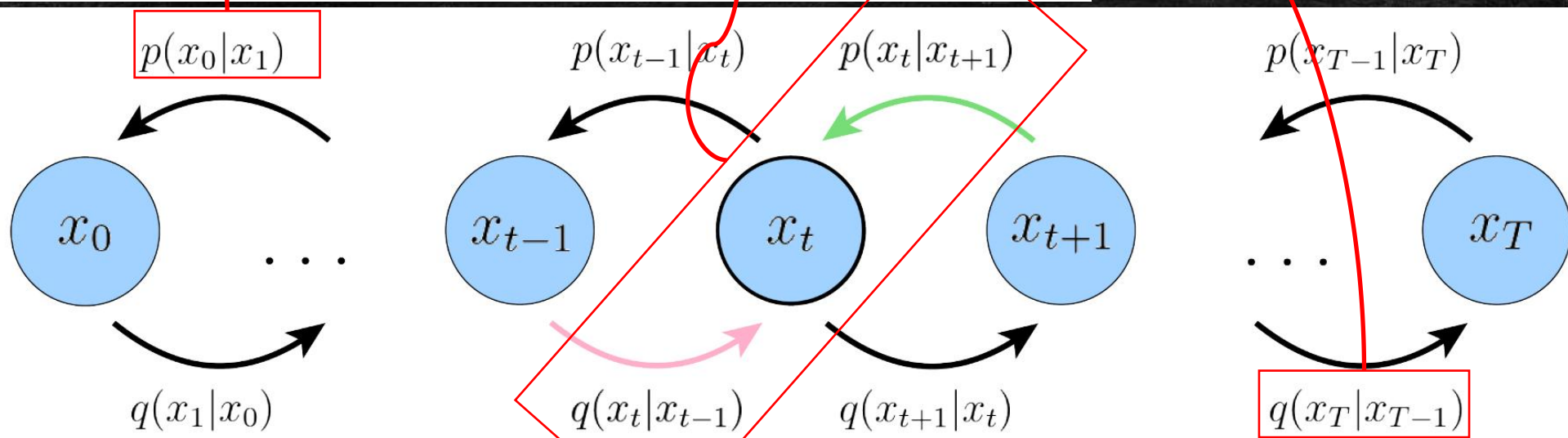
$$\begin{aligned} & \mathbb{E}_{q(x_{t-1}, x_t, x_{t+1}|x_0)} \left[\log \frac{p_\theta(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] \quad \boxed{\text{Consistency term}} \\ &= \iiint \left[\log \frac{p_\theta(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] q(x_{t-1}, x_t, x_{t+1}|x_0) dx_{t-1} dx_t dx_{t+1} \quad \rightarrow \boxed{\text{多元函数的期望定义}} \\ &= \iiint \left[\log \frac{p_\theta(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] q(x_{t-1}, x_{t+1}|x_0) dx_{t-1} dx_{t+1} q(x_t|x_{t-1}, x_{t+1}, x_0) dx_t \quad \rightarrow \boxed{\text{链式法则}} \\ &= \iint \left[\mathbb{E}_{q(x_t|x_{t-1})} \log \frac{p_\theta(x_t|x_{t+1})}{q(x_t|x_{t-1})} \right] q(x_{t-1}, x_{t+1}|x_0) dx_{t-1} dx_{t+1} \quad \rightarrow \boxed{\text{期望的定义+马尔科夫性质}} \\ &= - \iint [D_{KL}(q(x_t|x_{t-1}) || p_\theta(x_t|x_{t+1}))] q(x_{t-1}, x_{t+1}|x_0) dx_{t-1} dx_{t+1} \quad \rightarrow \boxed{\text{KL散度的定义}} \\ &= - \mathbb{E}_{q(x_{t-1}, x_{t+1}|x_0)} [D_{KL}(q(x_t|x_{t-1}) || p_\theta(x_t|x_{t+1}))] \quad \rightarrow \boxed{\text{期望的定义}} \end{aligned}$$

VDM数学推导

VDM的ELBO拆解，看其具体的含义

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))]}_{\text{prior matching term}} \\ &\quad - \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}}\end{aligned}$$

对比VAE:
都有Reconstruction term和
prior matching term
新增consistency term, 且该项
占主导, 因为包含很多子项



VDM数学推导

VDM的ELBO拆解，看其具体的含义

$$- \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1} | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t | \mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))]}_{\text{consistency term}}$$

包含 \mathbf{x}_{t-1} 和 \mathbf{x}_{t+1} 两个随机变量去
预估 \mathbf{x}_t
用蒙特卡洛方法计算这一项期望
误差，两个随机变量估计会方差
较大

有办法把2个随机变量变成1个吗？

可以！

通过如下的贝叶斯公式：

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t | \mathbf{x}_0)}{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}$$

VDM数学推导

2个随机变量变成1个随机变量推导(1/2页):

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] && \rightarrow \text{ELBO} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] && \rightarrow \text{链式法则+马尔科夫性质} \quad \text{之前推导过} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] && \rightarrow \text{分子分母连乘里各拆出一项} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] && \rightarrow \text{马尔科夫性质} \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\boxed{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)}} \right] && \rightarrow \log(AB) = \log(A) + \log(B) \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\boxed{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}}} \right] && \rightarrow \text{贝叶斯公式}\end{aligned}$$

VDM数学推导

2个随机变量变成1个随机变量推导(2/2页):

$$\begin{aligned}
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{\cancel{q(\mathbf{x}_1|\mathbf{x}_0)}} + \log \frac{\cancel{q(\mathbf{x}_1|\mathbf{x}_0)}}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \rightarrow \boxed{\log(AB) = \log(A) + \log(B)} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \rightarrow \boxed{\log(A) + \log(B) = \log(AB)} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \rightarrow \boxed{E[\log(AB)] = E[\log(A)] + E[\log(B)]} \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \rightarrow \boxed{\text{删掉无关的变量}} \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \rightarrow \boxed{\text{KL散度的定义}}
 \end{aligned}$$

VDM数学推导

VDM的ELBO拆解, Denoising matching term最后一步推导:

$$\begin{aligned} & \mathbb{E}_{q(x_t, x_{t-1}|x_0)} \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] \quad \boxed{\text{Denoising matching term}} \quad \begin{array}{l} \text{推导和Prior matching term} \\ \text{最后一步推导类似} \end{array} \\ &= \iint \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] q(x_t, x_{t-1}|x_0) dx_{t-1} dx_t \quad \rightarrow \quad \boxed{\text{多元函数的期望定义}} \\ &= \iint \left[\log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] q(x_{t-1}|x_t, x_0) dx_{t-1} q(x_t|x_0) dx_t \quad \rightarrow \quad \boxed{\text{链式法则}} \\ &= \int \left[\mathbb{E}_{q(x_{t-1}|x_t, x_0)} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right] q(x_t|x_0) dx_t \quad \rightarrow \quad \boxed{\text{期望的定义}} \\ &= - \int [D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))] q(x_t|x_0) dx_t \quad \rightarrow \quad \boxed{\text{KL散度的定义}} \\ &= - \mathbb{E}_{q(x_t|x_0)} [D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))] \quad \rightarrow \quad \boxed{\text{期望的定义}} \end{aligned}$$

VDM数学推导

VDM的ELBO拆解， 2个随机变量变成1个随机变量

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}\end{aligned}$$

对比之前2个随机变量公式：

- Reconstruction term和prior matching term保持一致
- 原来的consistency term，变成了denoising matching term
- 原来consistency term由 x_{t-1} 和 x_{t+1} 2个随机变量去预估 x_t
- 现在denoising matching term仅由 x_{t+1} 1个随机变量去预估 x_t ，采用蒙特卡洛方法时，估计方差变小

说明：

以上公式仅仅用到了马尔科夫定理，所以适用于所有MHVAE模型，而不仅仅局限于VDM模型。

VDM数学推导

VDM的ELBO拆解， 2个随机变量变成1个随机变量

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}\end{aligned}$$

对比VAE推导的ELBO拆解公式：

$$\mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right] = \underbrace{\mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \parallel p(z))}_{\text{prior matching term}}$$

结论：

当 $T=1$ 时，denosing matching term=0， x_0 用 x 表示， x_1 用 z 表示，VDM的ELBO拆解与VAE的ELBO拆解结果一致。

VDM数学推导

VDM的ELBO拆解, 1个随机变量版本

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

无可优化参数

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

denoising matching term由于是多项的求和, 所以会在优化目标中占主要部分。所以重点看看这部分吧!

$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ 这个分布是模型要去学习的。

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ 这个分布就是我要让模型学习的目标, 尽可能与这个分布相等。换句话说: **$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ 就是模型要学习的ground-truth。**

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ 这个GT怎么算呢?

VDM数学推导

$q(x_{t-1}|x_t, x_0)$ 这个ground-truth怎么算呢?

Recall: MHVAE \rightarrow VDM的三个限制条件之二

所有的encoder $q(x_t|x_{t-1})$ 都不需要学习, 而是预定义好的高斯分布模型, 就是 x_t 为以 x_{t-1} 为均值的高斯分布, 用数学表达就是:

$$q(x_t|x_{t-1}) \sim N(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

看到 $q(x_t|x_{t-1})$, 求 $q(x_{t-1}|x_t, x_0)$, **贝叶斯公式!**

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

另外几项: $q(x_t|x_{t-1}, x_0)$, $q(x_{t-1}|x_0)$ 和 $q(x_t|x_0)$ 具体怎么计算呢?

再**Recall: 重参数化技巧**

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon \quad \text{with } \epsilon \sim N(\epsilon; 0, I)$$

VDM数学推导

$q(x_{t-1}|x_t, x_0)$ 这个ground-truth怎么算呢?

用递推算 $q(x_t|x_0)$:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$$

重参数化

$$= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2}^* \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^* + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \epsilon_{t-2}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \epsilon_{t-2}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-2}$$

$= \dots$

$$= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0$$

$$= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$$

$$\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

x_{t-1} 重参数化+迭代

系数合并

高斯分布: 和的方差等于方差的和
和的均值等于均值的和

系数合并

x_{t-2}, \dots, x_1 重参数化+多次迭代直到用 x_0 来表示

参数表示替换: $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

$q(x_t|x_0)$ 为高斯分布

VDM数学推导

$q(x_{t-1}|x_t, x_0)$ 这个ground-truth怎么算呢? (推导1/2页)

用同样的方法也能得到: $q(x_{t-1}|x_0) \sim N(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})I)$

$q(x_t|x_{t-1}, x_0)$, $q(x_{t-1}|x_0)$ 和 $q(x_t|x_0)$ 三项带入贝叶斯公式:

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)}$$



贝叶斯公式

$$= \frac{\mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I) \mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})I)}{\mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)}$$



带入三个高斯分布

$$\propto \exp \left\{ - \left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{2(1 - \alpha_t)} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{2(1 - \bar{\alpha}_{t-1})} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{2(1 - \bar{\alpha}_t)} \right] \right\}$$



带入高斯分布的表达式, 隐去系数

$$= \exp \left\{ - \frac{1}{2} \left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{1 - \alpha_t} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)^2}{1 - \bar{\alpha}_t} \right] \right\}$$



提出通项1/2

$$= \exp \left\{ - \frac{1}{2} \left[\frac{(-2\sqrt{\alpha_t}x_t x_{t-1} + \alpha_t x_{t-1}^2)}{1 - \alpha_t} + \frac{(x_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}x_{t-1}x_0)}{1 - \bar{\alpha}_{t-1}} + C(x_t, x_0) \right] \right\}$$



平方开方

$$\propto \exp \left\{ - \frac{1}{2} \left[- \frac{2\sqrt{\alpha_t}x_t x_{t-1}}{1 - \alpha_t} + \frac{\alpha_t x_{t-1}^2}{1 - \alpha_t} + \frac{x_{t-1}^2}{1 - \bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}x_{t-1}x_0}{1 - \bar{\alpha}_{t-1}} \right] \right\}$$



先忽略常数

$$= \exp \left\{ - \frac{1}{2} \left[\left(\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\}$$



合并同类项系数

VDM数学推导

$q(x_{t-1}|x_t, x_0)$ 这个ground-truth怎么算呢? (推导2/2页)

$$\begin{aligned}
 &= \exp \left\{ -\frac{1}{2} \left[\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t} x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} x_0}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} && \rightarrow \text{合并同类项系数} \\
 &= \exp \left\{ -\frac{1}{2} \left(\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t} x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} x_0}{1 - \bar{\alpha}_{t-1}} \right)}{\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}} x_{t-1} \right] \right\} && \rightarrow \text{提取系数} \\
 &= \exp \left\{ -\frac{1}{2} \left(\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t} x_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} x_0}{1 - \bar{\alpha}_{t-1}} \right) (1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_{t-1} \right] \right\} && \rightarrow \text{系数化简} \\
 &= \exp \left\{ -\frac{1}{2} \left(\frac{1}{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}} \right) \left[x_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} x_{t-1} \right] \right\} && \rightarrow \text{凑成高斯分布表达式的系数形式} \\
 &\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)}) && \rightarrow \text{把之前忽略的常数拿回来, 正好能配成一个高斯分布的表达式, 得到分布的均值和方差}
 \end{aligned}$$

高斯分布: $x \sim N(x; \mu, \sigma^2) =$

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2}\right) [x^2 - 2\mu x + \mu^2]\right) \propto \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2}\right) [x^2 - 2\mu x]\right)$$

VDM数学推导

$q(x_{t-1}|x_t, x_0)$ 这个ground-truth的特性?

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)})$$

- 高斯分布
- 方差 $\Sigma_q(t)$ 只跟的 α_t 有关, 每一步的 α_t 已知
- 均值 $\mu_q(x_t, x_0)$ 是关于 x_t, x_0 的函数

Recall: 求这个是干啥啊? 因为如下ELBO拆解的第三项需要用模型拟合 $q(x_{t-1}|x_t, x_0)$

$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

模型优化目标: 最小化 $p_{\theta}(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度

VDM数学推导

最小化 $p_\theta(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)})$$

- 高斯分布
- 方差 $\Sigma_q(t)$ 只跟的 α_t 有关, 每一步的 α_t 已知
- 均值 $\mu_q(x_t, x_0)$ 是关于 x_t, x_0 的函数

模型优化目标: 最小化 $p_\theta(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度

根据ground-truth的前两条特性, 我们也可以给我们预估值以下两条性质:

- ✓ 定义 $p_\theta(x_{t-1}|x_t)$ 为高斯分布
- ✓ 定义方差与 $q(x_{t-1}|x_t, x_0)$ 的相等, 等于 $\Sigma_q(t)$

VDM数学推导

最小化 $p_{\theta}(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度

下标 θ 和 q 分别代表模型预估的分布和ground-truth分布

$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))$$

$$= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \rightarrow \text{方差相等, 均为}\boldsymbol{\Sigma}_q(t)$$

$$= \arg \min_{\theta} \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_q(t)|}{|\boldsymbol{\Sigma}_q(t)|} - d + \text{tr}(\boldsymbol{\Sigma}_q(t)^{-1} \boldsymbol{\Sigma}_q(t)) + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \rightarrow \text{两个高斯分布的KL散度计算公式}$$

$$= \arg \min_{\theta} \frac{1}{2} [\log 1 - d + d + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)] \rightarrow \text{tr}(I) = d, d \text{ 是方差矩阵的维度}$$

$$= \arg \min_{\theta} \frac{1}{2} [(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)] \rightarrow \text{简单化简}$$

$$= \arg \min_{\theta} \frac{1}{2} [(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T (\sigma_q^2(t) \mathbf{I})^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)] \rightarrow \boldsymbol{\Sigma}_q(t) = \sigma_t^2 I \text{ 已知}$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} [\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2] \rightarrow \text{简单化简}$$

VDM数学推导

最小化 $p_{\theta}(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度

$$\arg \min_{\theta} D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} [\|\mu_{\theta} - \mu_q\|_2^2]$$

最小化KL散度 \Leftrightarrow 拟合 $q(x_{t-1}|x_t, x_0)$ 均值

要拟合的均值 $\mu_q(x_t, x_0)$ 如下:

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

$\mu_{\theta}(x_t, t)$ 也是关于 x_t 的表达式, 为了接近 $\mu_q(x_t, x_0)$ 的形式, 我们将 $\mu_{\theta}(x_t, t)$ 写成如下公式:

$$\mu_{\theta}(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{x}_{\theta}(x_t, t)}{1 - \bar{\alpha}_t}$$

这样: 最小化KL散度 \Leftrightarrow 拟合 $q(x_{t-1}|x_t, x_0)$ 均值 \Leftrightarrow 根据 x_t 和步数 t 拟合 x_0

VDM数学推导

最小化 $p_{\theta}(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度

最小化KL散度 \Leftrightarrow 拟合 $q(x_{t-1}|x_t, x_0)$ 均值 \Leftrightarrow 根据 x_t 和步数 t 拟合 x_0

$$\begin{aligned} \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)) &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} - \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \right\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} - \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \right\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} (\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0) \right\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} \left[\|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2 \right] \end{aligned}$$

全部都是简单化简

VDM数学推导

最小化 $p_{\theta}(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度
(预测噪声的损失函数表达)

$$\arg \min_{\theta} D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} [\|\mu_{\theta} - \mu_q\|_2^2]$$

Recall:

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} \Rightarrow x_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_0}{\sqrt{\bar{\alpha}_t}}$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_0}{\sqrt{\bar{\alpha}_t}}}{1 - \bar{\alpha}_t}$$

简单带入

$$= \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0$$

合并同类项 & $\bar{\alpha}_t = \bar{\alpha}_{t-1} \alpha_t$

同理：我们把模型预测的 $\mu_{\theta}(x_t, t)$ 也写成上面的形式：

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_{\theta}(x_t, t)$$

VDM数学推导

最小化 $p_{\theta}(x_{t-1}|x_t)$ 和 $q(x_{t-1}|x_t, x_0)$ 的KL散度
(关于预测噪声的损失函数表达)

$$\begin{aligned}\arg \min_{\theta} D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t)) &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_{\theta} - \mu_q\|_2^2 \right] \\&= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_{\theta}(x_t, t) - \frac{1}{\sqrt{\alpha_t}} x_t + \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0 \right\|_2^2 \right] \\&= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0 - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_{\theta}(x_t, t) \right\|_2^2 \right] \\&= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} (\epsilon_0 - \hat{\epsilon}_{\theta}(x_t, t)) \right\|_2^2 \right] \\&= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t) \alpha_t} \left[\|\epsilon_0 - \hat{\epsilon}_{\theta}(x_t, t)\|_2^2 \right]\end{aligned}$$

全部都是简单化简

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0$$

这样: 最小化KL散度 \Leftrightarrow 拟合 $q(x_{t-1}|x_t, x_0)$ 均值 \Leftrightarrow 拟合初始图像 x_0 \Leftrightarrow 拟合噪声 ϵ_0

VDM数学推导

Inference阶段:

每次来一个 x_t 和 t , 模型预测得到加的噪声: $\hat{\epsilon}_\theta(x_t, t)$

然后就可以算 x_{t-1} 的均值, 然后标准高斯分布上随机采样出一个噪声 \mathbf{z} , 方差 $\Sigma_q(t) = \sigma_t^2 I$ 已知, 用重参数化技巧得到 x_{t-1} :

$$x_{t-1}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_\theta(x_t, t) + \sigma_t \mathbf{z}$$

经过 t 次迭代, 就能得到没有噪声的 x_0

DDPM训练和采样流程:

Algorithm 1 Training

- 1: **repeat**
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on
 $\text{系数} \nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$
- 6: **until** converged

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ **if** $t > 1$, **else** $\mathbf{z} = \mathbf{0}$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return** \mathbf{x}_0

小结

- 1、VDM = 马尔可夫级联VAE (MHVAE) + 3个限制条件
- 2、优化VDM和优化VAE一样，都是最大化ELBO
- 3、拆解VDM的ELBO能得到Reconstruction term、prior matching term和denoising matching term三项，第三项为VDM相比VAE新增，该项因为是多项求和，占据损失函数的主导。当 $T=1$ 时，VDM和VAE的ELBO相等。

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \end{aligned}$$

- 4、 $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ 就是模型要学习的ground-truth。这一项可以计算得到。也是高斯分布，满足：

$$\mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)})$$

小结

5、模型优化目标：最小化KL散度 \Leftrightarrow 拟合 $q(x_{t-1}|x_t, x_0)$ 均值 \Leftrightarrow 拟合初始图像 $x_0 \Leftrightarrow$ 拟合噪声 ϵ_0

$$\begin{aligned} \arg \min_{\theta} D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t)) &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_{\theta} - \mu_q\|_2^2 \right] \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} \left[\|\hat{x}_{\theta}(x_t, t) - x_0\|_2^2 \right] = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \left[\|\epsilon_0 - \hat{\epsilon}_{\theta}(x_t, t)\|_2^2 \right] \end{aligned}$$

6、Inference时，根据模型预测计算t-1步均值，再在标准高斯分布上随机采样出一个噪声 z ，方差 $\Sigma_q(t) = \sigma_t^2 I$ 已知，用重参数化技巧得到 x_{t-1} ：

$$\begin{aligned} \mu_{\theta}(x_t, t) &= \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_{\theta}(x_t, t) \\ x_{t-1}(x_t, t) &= \mu_{\theta}(x_t, t) + \sigma_t z = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_{\theta}(x_t, t) + \sigma_t z \end{aligned}$$

结语

基础理论大厦已经建成，但是仅凭如此还不够！

- 迭代步数太多，生成速度太慢 (DDIM)
- 隐空间维度低，无法生成高清图 (SD使用的LDM)
- 加噪过程能否优化，固定方差，只预测均值真的是最优解吗 (IDDPM)
- 如何控制生成的方向，得到想要的生成内容 (CG、CFG)
- 什么网络结构能得到更好的生成效果 (U-net、DiTs)
- 更多实验得到的工程优化 (DDPM的L-simple)
-