

一、条件概率公式与高斯分布的KL散度

1、条件概率的一般形式

$$P(A, B, C) = P(C \mid B, A)P(B, A) = P(C \mid B, A)P(B \mid A)P(A)$$

$$P(B, C \mid A) = P(B \mid A)P(C \mid A, B)$$

2、基于马尔科夫假设的条件概率

如果满足马尔科夫链关系A->B->C, 那么有

$$P(A, B, C) = P(C \mid B, A)P(B, A) = P(C \mid B)P(B \mid A)P(A)$$

$$P(B, C \mid A) = P(B \mid A)P(C \mid B)$$

3、高斯分布的KL散度公式

对于两个单一变量的高斯分布p和q而言，它们的KL散度为

$$KL(p, q) = \log \frac{\sigma_1^2}{\sigma_2^2} + \frac{\sigma^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

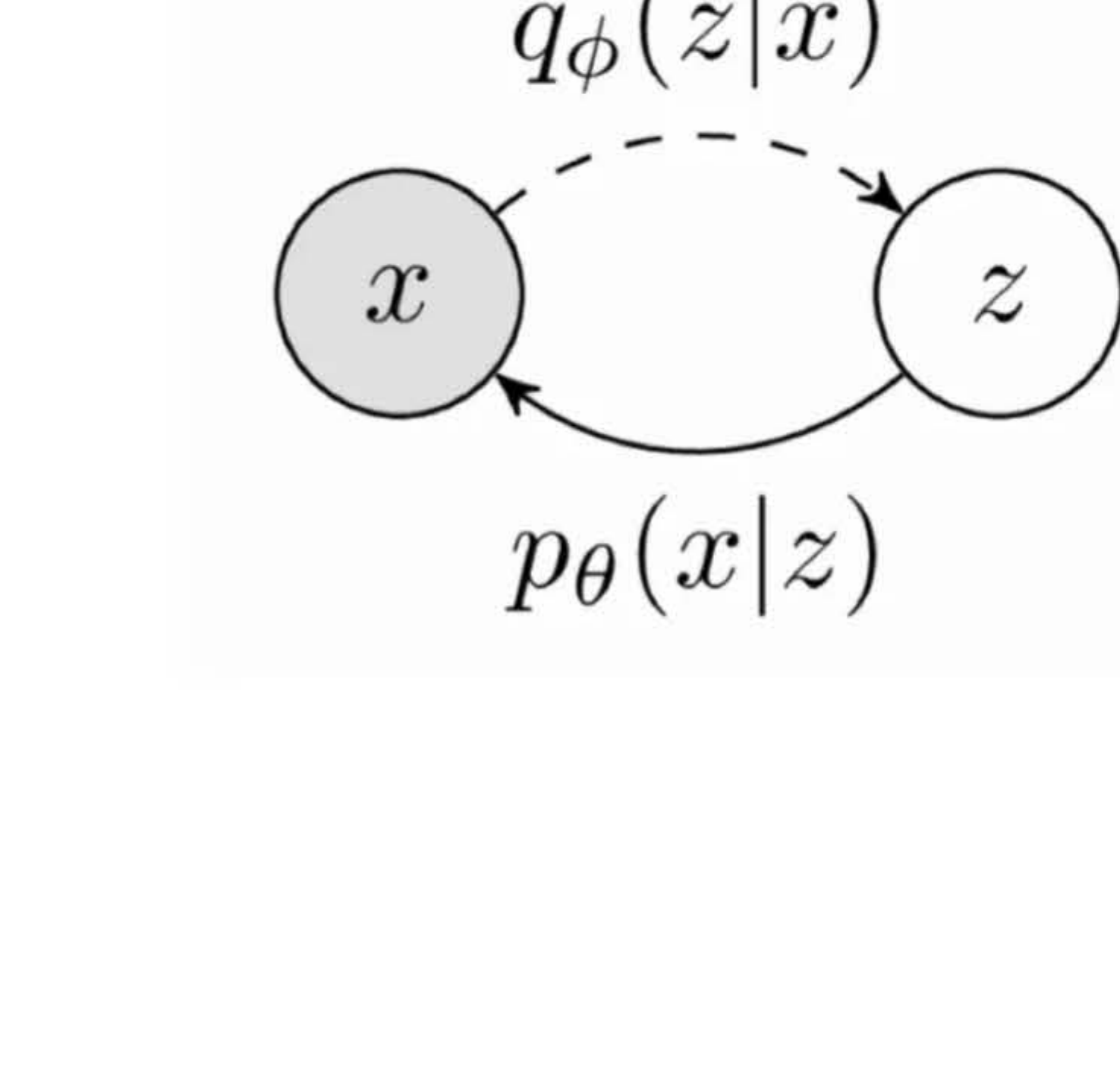
4、重参数技巧

若希望从高斯分布N(μ, σ)中采样，可以先从标准分布N(0, 1)采样出z，再得到σ * z + μ。这样做的好处是将随机性转移到了z这个常量上，而σ和μ则当做仿射变换网络的一部分。



二、VAE与多层VAE回顾

1、单层VAE的原理公式与置信下界



q加噪, p去噪
 $q_\phi(z|x)$ 表示函数名q, 给定x条件下z的分布函数.
分布函数的参数是 ϕ
所以机器学习里面的函数都可以看做 $f_\theta(y|x)$
这种分布.

$$p(x) = \int_z p_\theta(x|z)p(z)$$

因为 $p_\theta(x|z) = \frac{p_\theta(x, z)}{p(z)}$ 所以这个式子是边缘分布的积分.

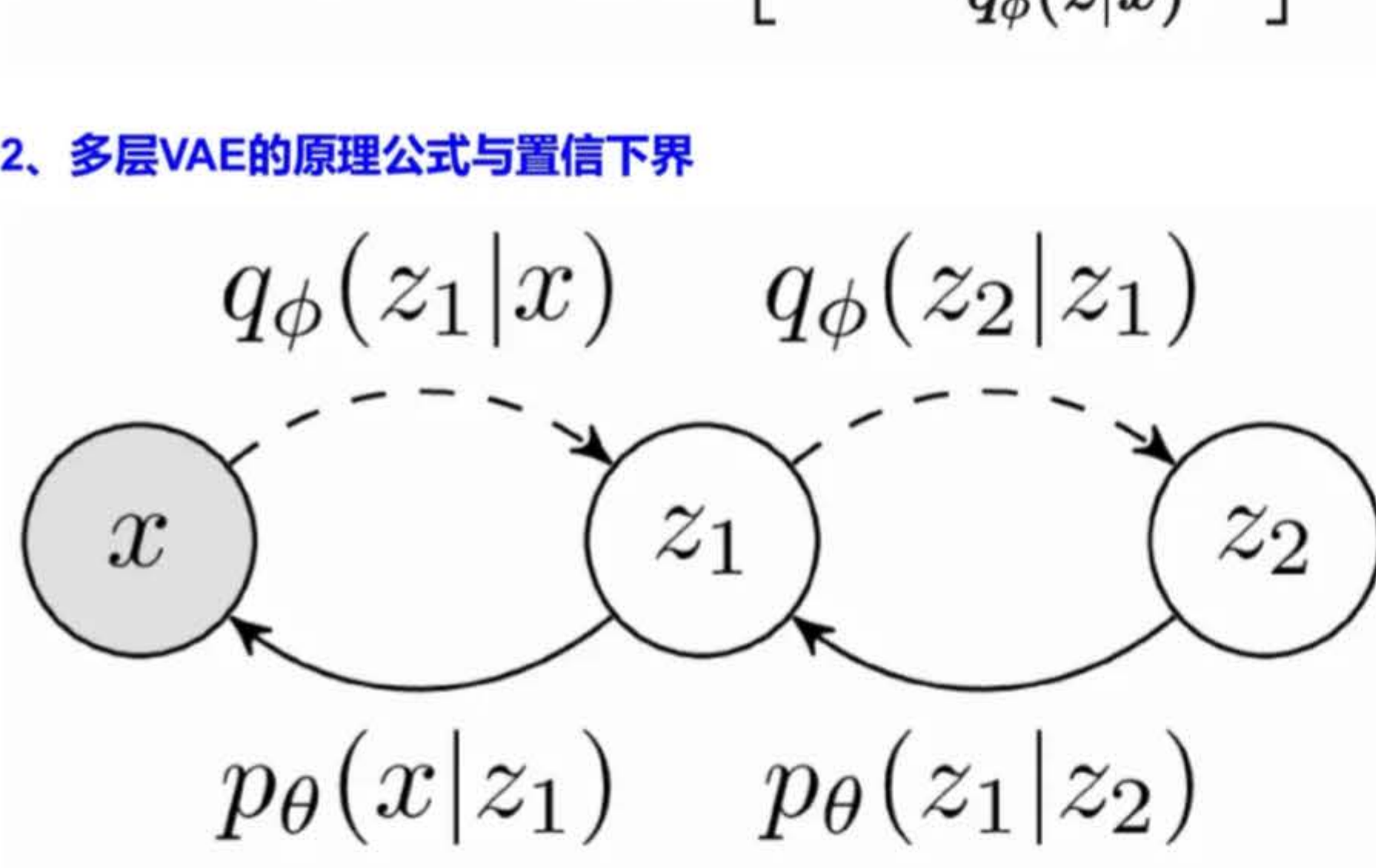
$$p(x) = \int q_\phi(z|x) \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)}$$

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

积分里面的z分布函数就可以提出来当期望函数.

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

2、多层VAE的原理公式与置信下界



$$p(x) = \int_{z_1} \int_{z_2} p_\theta(x, z_1, z_2) dz_1 dz_2$$

$$p(x) = \int \int q_\phi(z_1, z_2|x) \frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)}$$

$$p(x) = \mathbb{E}_{z_1, z_2 \sim q_\phi(z_1, z_2|x)} \left[\frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)} \right]$$

$$\log p(x) \geq \mathbb{E}_{z_1, z_2 \sim q_\phi(z_1, z_2|x)} \left[\log \frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)} \right]$$

$$p(x, z_1, z_2) = p(x|z_1)p(z_1|z_2)p(z_2)$$

$$q(z_1, z_2|x) = q(z_1|x)q(z_2|z_1)$$

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(z_1, z_2|x)} [\log p(x|z_1) - \log q(z_1|x) + \log p(z_2|z_1) - \log q(z_2|z_1) + \log p(z_2)]$$

三、Diffusion Model图示

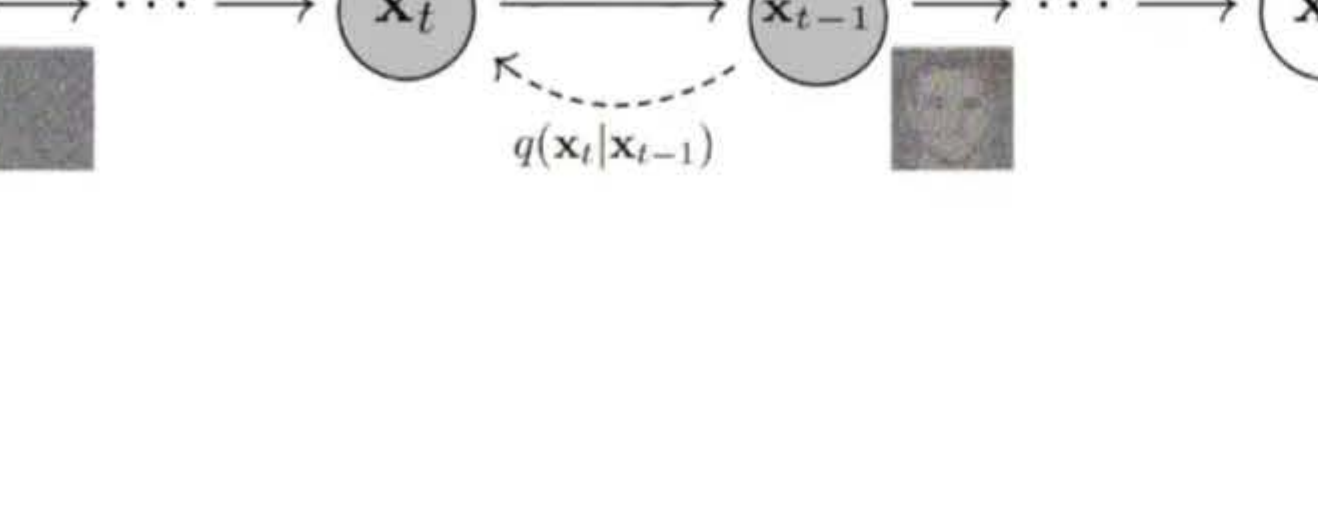
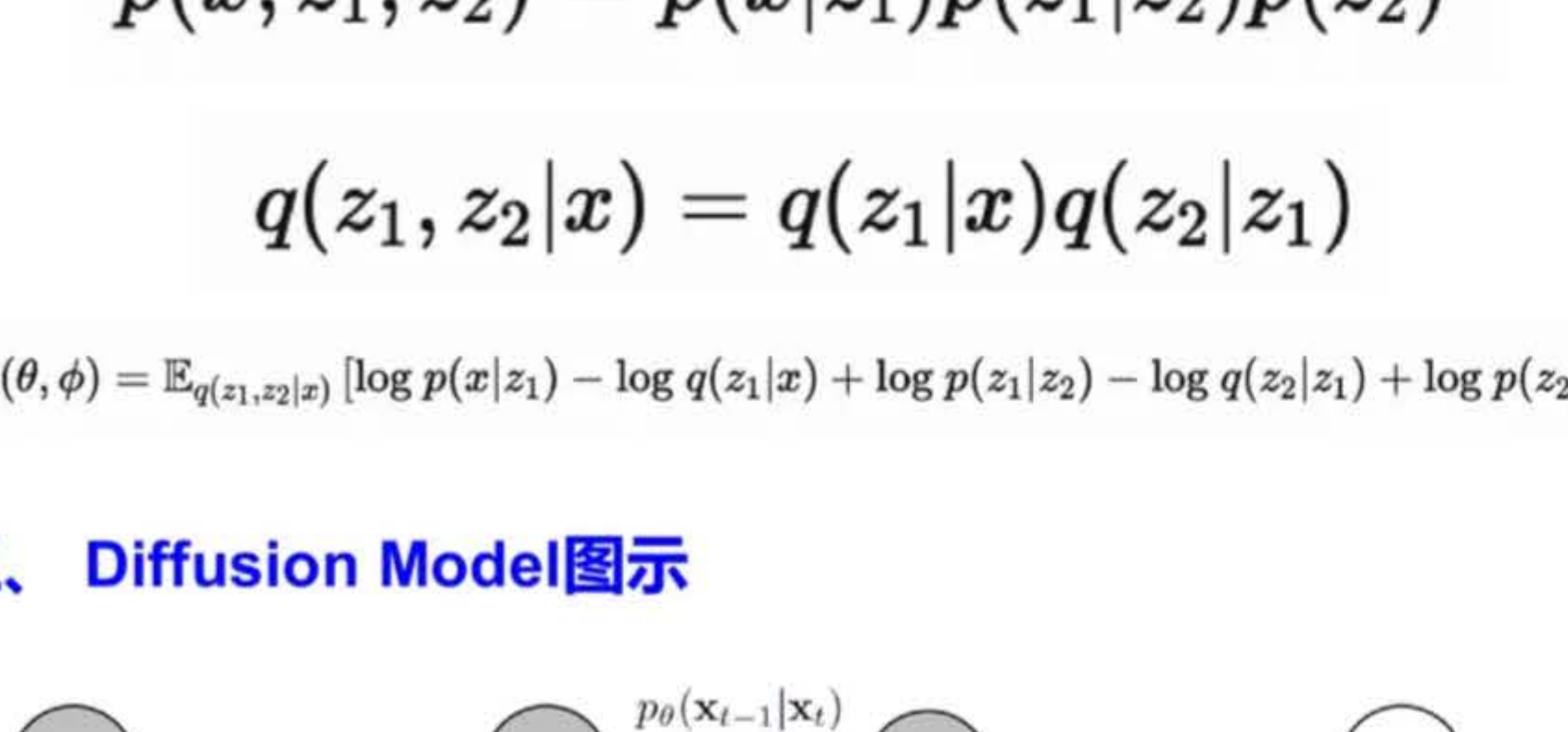


Figure 1. The proposed modeling framework trained on 2-d swiss roll data. The top row shows time slices from the forward trajectory $q(x^{(0:T)})$. The data distribution (left) undergoes Gaussian diffusion, which gradually transforms it into an identity-covariance Gaussian (right). The middle row shows the corresponding time slices from the trained reverse trajectory $p(x^{(0:T)})$. An identity-covariance Gaussian (right) undergoes a Gaussian diffusion process with learned means and covariance functions, and is gradually transformed back into the data distribution (left). The bottom row shows the drift term, $f_{\theta}(x^{(t)}, t) - x^{(t)}$, for the same reverse diffusion process.

四、扩散过程 (Diffusion Process)

1、给定初始数据分布 $x_0 \sim q(x)$ ，可以不断地向分布中添加高斯噪声，该噪声的标准差是以固定值 β_t 而确定的，均值是以固定值 β_t 和当前时刻的数据 x_t 决定的。这个过程是一个马尔科夫链过程。 → 最终目标为了刻画 $q(x)$

2、随着t的不断增大，最终数据分布 x_T 变成了一个各向独立的高斯分布。

Given a data point sampled from a real data distribution $x_0 \sim q(x)$, let us define a forward diffusion process in which we add small amount of Gaussian noise to the sample in T steps, producing a sequence of noisy samples x_1, \dots, x_T . The step sizes are controlled by a variance schedule $\{\beta_t\}_{t=1}^T$.

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t) \quad q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

The data sample x_0 gradually loses its distinguishable features as the step t becomes larger. Eventually when $T \rightarrow \infty$, x_T is equivalent to an isotropic Gaussian distribution.

3、任意时刻的 $q(x_t)$ 推导也可以完全基于 x_0 和 β_t 来计算出来，而不需要做迭代

注意这里，两个正态分布 $X \sim \mathcal{N}(\mu_1, \sigma_1)$ 和 $Y \sim \mathcal{N}(\mu_2, \sigma_2)$ 的叠加后的分布 $aX + bY$ 的均值为 $a\mu_1 + b\mu_2$ ，方差为 $a^2\sigma_1^2 + b^2\sigma_2^2$ ，所以 $\sqrt{\alpha_t - \alpha_{t-1}}z_{t-2} + \sqrt{1 - \alpha_t}z_{t-1}$ 可以重参数化成只含一个随机变量 z 构成的 $\sqrt{1 - \alpha_t}\alpha_{t-1}z$ 的形式。

A nice property of the above process is that we can sample x_t at any arbitrary time step t in a closed form using **reparameterization trick**. Let $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$:

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}z_t \quad ; \text{where } z_{t-1}, z_{t-2}, \dots \sim \mathcal{N}(0, \mathbf{I})$$

$$= \sqrt{\alpha_t\bar{\alpha}_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\bar{\alpha}_{t-1}}z_{t-2} \quad ; \text{where } z_{t-2} \text{ merges two Gaussians (*)}$$

$$= \dots$$

$$= \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}z$$

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_0, (1 - \alpha_t)\mathbf{I})$$

(*) Recall that when we merge two Gaussians with different variance, $\mathcal{N}(0, \sigma_1^2\mathbf{I})$ and $\mathcal{N}(0, \sigma_2^2\mathbf{I})$, the new distribution is $\mathcal{N}(0, (\sigma_1^2 + \sigma_2^2)\mathbf{I})$. Here the merged standard deviation is $\sqrt{(1 - \alpha_t) + \alpha_t(1 - \alpha_{t-1})} = \sqrt{1 - \alpha_t\alpha_{t-1}}$.

Usually, we can afford a larger update step when the sample gets noisier, so $\beta_1 < \beta_2 < \dots < \beta_T$ and therefore $\bar{\alpha}_1 > \dots > \bar{\alpha}_T$. → 因为beta1距离目标越远所以越需要降低学习难度.

五、逆扩散过程 (Reverse Process)

逆过程是从高斯噪声中恢复原始数据，我们可以假设它也是一个高斯分布，但没法直接去拟合分布，所以需要构建一个参数分布去做估计，逆扩散过程仍然是一个马尔科夫链过程。

If we can reverse the above process and sample from $q(x_{t-1}|x_t)$, we will be able to recreate the true sample from a Gaussian noise input, $x_T \sim \mathcal{N}(0, \mathbf{I})$. Note that if β_t is small enough, $q(x_{t-1}|x_t)$ will also be Gaussian. Unfortunately, we cannot easily estimate $q(x_{t-1}|x_t)$ because it needs to use the entire dataset and therefore we need to learn a model p_θ to approximate these conditional probabilities in order to run the reverse diffusion process.

$$p_\theta(x_0:T) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \quad p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

六、后验的扩散条件概率 $q(x_{t-1}|x_t, x_0)$ 分布是可以由公式表达的

多加一个 x_0 的话有解析式.

也就是说，给定 x_t 和 x_0 ，我们可以计算出 x_{t-1}

注意：高斯分布的概率密度函数是 $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

注意： $ax^2 + bx = a(x + \frac{b}{2a})^2 + C$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \hat{\mu}(x_t, x_0), \hat{\beta}_t\mathbf{I})$$

Using Bayes' rule, we have:

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

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

上面公式带入得到这个.

where $C(x_t, x_0)$ is some function not involving x_{t-1} and details are omitted. Following the standard Gaussian density function, the mean and variance can be parameterized as follows:

$$\hat{\beta}_t = \frac{1}{\beta_t} \left(\frac{\alpha_t}{1 - \bar{\alpha}_{t-1}} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) = \frac{1 - \bar{\alpha}_t}{1 - \bar{\alpha}_{t-1}} \beta_t$$

$$\hat{\mu}(x_t, x_0) = \left(\frac{\sqrt{\alpha_t}x_t + \sqrt{1 - \alpha_t}x_0 \right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\alpha_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0$$

变换成一个高斯函数.

!!!紧急重要更新：大家注意，上面蓝色公式的第二项 x_0 系数中的分子分母中的所有 $\bar{\alpha}_t$ 应该写成 $\bar{\alpha}_{t-1}$ 。

估计是因为 α_t 与 α_{t-1} 差别很小.所以这里面都记做 α_t 了.

根据前面 x_0 与 x_t 之间的关系式，我们可以知道

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}z_t)$$

将 x_0 的表达式代入到 $q(x_{t-1}|x_t, x_0)$ 的分布中，可以重新给出该分布的均值表达式，这个时候表达式中不再含有 x_0 ，并且多了噪声项，这为我们设计神经网络提供了基础。也就是说，在给定 x_0 的条件下，后验

条件高斯分布的均值计算只与 x_t 和 z_t 有关。 z_t 是t时刻的随机正态分布变量，源自重参数化。

$$\hat{\mu}_t = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\alpha_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \frac{1}{\sqrt{\alpha_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}z_t)$$

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} z_t \right)$$

七、目标数据分布的似然函数

我们可以在负对数似然函数的基础上加上一个KL散度，于是就构成了负对数似然的升级版了，升级版，负对数似然自然也就越小，那么对数似然就越大了。

$$-\log p_\theta(x_0) \leq -\log p_\theta(x_0) + D_{KL}(q(x_{1:T}|x_0) \| p_\theta(x_{1:T}|x_0)) \rightarrow \text{因为KL散度横大于0}$$

$$= -\log p_\theta(x_0) + \mathbb{E}_{x_{1:T} \sim q(x_{1:T}|x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right]$$

$$= -\log p_\theta(x_0) + \mathbb{E}_q \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} + \log p_\theta(x_0) \right]$$

$$= \mathbb{E}_q \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right]$$

Let $L_{VLB} = \mathbb{E}_q(x_{0:T}) \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right] \geq -\mathbb{E}_q(x_0) \log p_\theta(x_0)$

进一步可以写出如上公式的交叉熵的上界，接下来，我们可以对交叉熵的上界进行化简

$$L_{VLB} = \mathbb{E}_{q(x_0)} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right]$$

$$= \mathbb{E}_q \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_T) \prod_{t=1}^{T-1} p_\theta(x_{t-1}|x_t)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=1}^T \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \sum_{t=2}^T \log \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \log \frac{q(x_T|x_0)}{q(x_1|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

$$= \mathbb{E}_q \left[\log \frac{q(x_T|x_0)}{p_\theta(x_T)} + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} - \log p_\theta(x_0|x_1) \right]$$

$$= \underbrace{\mathbb{E}_q \left[D_{KL}(q(x_T|x_0) \| p_\theta(x_T)) \right]}_{L_{\theta T}} + \sum_{t=2}^T \underbrace{D_{KL}(q(x_{t-1}|x_t, x_0) \| p_\theta(x_{t-1}|x_t)) - \log p_\theta(x_0|x_1)}_{L_{\theta t-1}} = L_\theta$$

这里论文将 $p_\theta(x_{t-1}|x_t)$ 分布的方差设置成一个与 β 相关的常数，因此可训练的参数只存在于其均值中

这个非常关键.

对于两个单一变量的高斯分布p和q而言，它们的KL散度为

$$KL(p, q) = \log \frac{\sigma_1^2}{\sigma_2^2} + \frac{\sigma^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

$$L_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \left\| \hat{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \right\|^2 \right] + C$$

既然这里的loss是从KL divergence出发的，或者说是与分布有关的，那我们可以设计一个黑箱子神经网络，把它称之为 D_θ 网络。对于 D_θ 网络，输入是 x_t 和时间编码 t ，对于输出是什么，取决于我们的建模目标。

而我们可以有多种建模目标：【1】直观的做法是让 D_θ 网络的输出等于前向过程中的后验分布均值 $\hat{\mu}(x_{t-1}|x_t, x_0)$ ，这种建模方法俗称预测后验分布的期望值；【2】根据 $\hat{\mu}(x_{t-1}|x_t, x_0)$ 表达式，它里面的 x_0 对于 D_θ 网络是未知的，因此第二种做法是让 D_θ 网络的输出等于 x_0 ，这种做法即直接预测原始数据。有人问，既然可以通过 D_θ 网络直接预测 x_0 了，那是不是采样过程就直接让 $D_\theta(x_T, T)$ 的输出即认为是生成了样本了呢？答案是直接一步到位，质量会比较差，还是需要通过马尔科夫高斯条件迭代而获得最终高质量的生产样本；【3】当我们把 $\hat{\mu}(x_{t-1}|x_t, x_0)$ 中的 x_0 用 x_t 去表示的时候， $\hat{\mu}(x_{t-1}|x_t, x_0)$ 就变成了如下只包含 x_t 和随机变量 z 的式子。其中 x_t 对于 D_θ 网络是已知的，而 z 是未知的，因此这个时候，我们可以选择建模目标是让 D_θ 网络的输出等于 z 了

$$L_{t-1} - C = \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2\sigma_t^2} \left\| \hat{\mu}_t \left(x_t(x_0, \epsilon), \frac{1}{\sqrt{\alpha_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon) \right) - \mu_\theta(x_t(x_0, \epsilon), t) \right\|^2 \right]$$

$$= \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t(x_0, \epsilon) - \sqrt{1 - \bar{\alpha}_t}\epsilon \right) - \mu_\theta(x_t(x_0, \epsilon), t) \right\|^2 \right]$$

在DDPM论文中，作者选择了方案【3】，即让 D_θ 网络的输出等于 z ，预测噪声均值。于是，新的逆向条件分布均值可以表示成(下式中的 ϵ_θ 相当于我们定义的广义的 D_θ 网络的具体分布形式)：

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

于是 L_{t-1} 可以化简成如下表达式

$$\mathbb{E}_{x_0, \epsilon} \left[2\sigma_t^2 \alpha_t \left(\frac{\beta_t^2}{1 - \bar{\alpha}_t} \right) \left\| \epsilon - \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2 \right]$$

DDPM作者又发现，干脆将系数丢掉，训练更加稳定质量更好，于是有了下面的 L_{simple}

$$L_{simple}(\theta) := \mathbb{E}_{t, x_0, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2 \right]$$