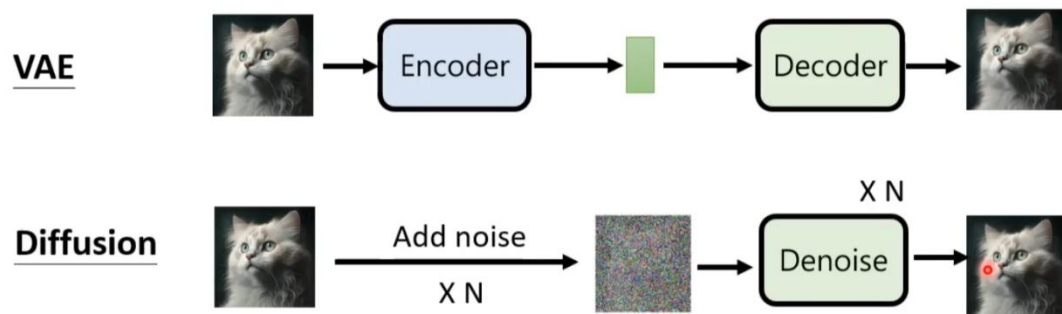


二. 模型背景与模型数学推导

我们刚刚已经介绍了图像生成这一问题，并且我们已经知道了扩散模型“图像生成”这一步骤是说的 **denoise** 过程。那他是怎么实现从一个随机的噪声中逐渐去噪变成想要的图像的呢？

VAE vs. Diffusion Model



主要演算法包含以下两个：

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
$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2$$
 - 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
-

他看起来非常非常的简单

Repeat 就是说一直做第二行到第五行的工作，直到 converge 为止



x_0 : clean image



ϵ : noise

我们需要首先在一个图像数据库中取样一个清晰的 x_0 图像，随后取样一个小 t ，小 t 是在 1 到大 T 中的一个整数。大 T 是一个比较大的数字，比如 1k

随后，我们在高斯分布中随机取样一个噪声，均值与方差是 0 和 1.

第五行算法相当复杂，我们一步一步拆分来看

第五行是说，我们需要把原图与噪声做一个加权，

Training

$\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_T$



x_0

$\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_T$

ϵ

Sample t

$$\sqrt{\bar{\alpha}_t} \begin{matrix} \text{cat image} \\ x_0 \end{matrix} + \sqrt{1 - \bar{\alpha}_t} \begin{matrix} \text{noise} \\ \epsilon \end{matrix} = \begin{matrix} \text{noisy cat image} \end{matrix}$$

这里的 α_1 到阿尔法 t 都是事先给好的，这样我们就得到了一个带噪声的图像，（这里的阿尔法 1 到阿尔法 t 是越来越小的），也就是说，给的 t 越大，噪声越凶狠，原图像约不清晰。得到噪声图像以后，我们发现他和 t 一起被放入了一个伊普西隆 θ ，这个就是 noise predictor

Algorithm 1 Training

1: repeat

2: $x_0 \sim q(x_0) \leftarrow \dots$ sample clean image

3: $t \sim \text{Uniform}(\{1, \dots, T\})$

4: $\epsilon \sim \mathcal{N}(0, I) \leftarrow \dots$ sample a noise

5: Take gradient descent step on

$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\underbrace{\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}_{\text{Noisy image}}, t) \right\|^2$$

6: until converged

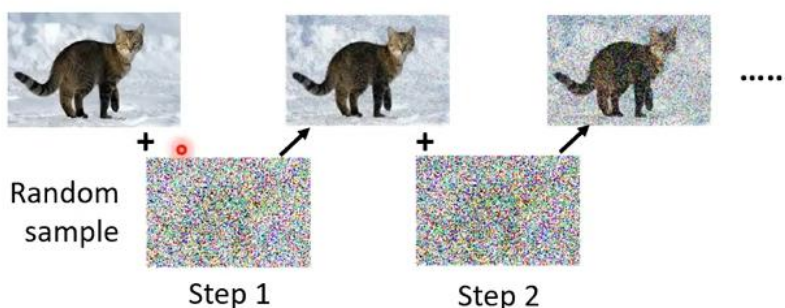
Noisy image

Noise
predictor

这个说的就是我预测噪声与原噪声的范数平方。梯度就是在这上面计算并训练的。也就是，给定噪声图像与 t ，我们相当于要去预测混入的噪音是怎样的。伊普西隆 θ 就

是用于预测这个噪音的。

想像中 ...



實際上 ...

$$\sqrt{\bar{\alpha}_t} x_n + \sqrt{1 - \bar{\alpha}_t} \varepsilon = \text{denoised image}$$

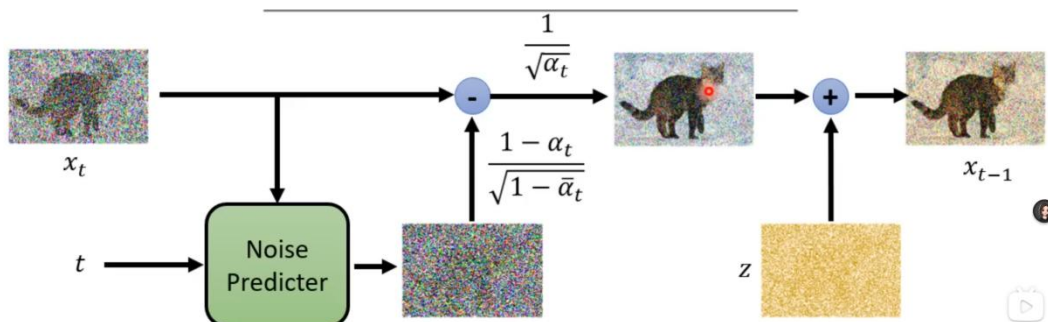
The diagram shows a clear cat image x_n multiplied by $\sqrt{\bar{\alpha}_t}$, plus a noise image ε multiplied by $\sqrt{1 - \bar{\alpha}_t}$, resulting in a denoised cat image.

但是按照我们第一部分讲的 这个噪声应该是一步一步加上的啊! 为什么训练过程是直接加上一个噪声呢!

这是我们的第一个疑问

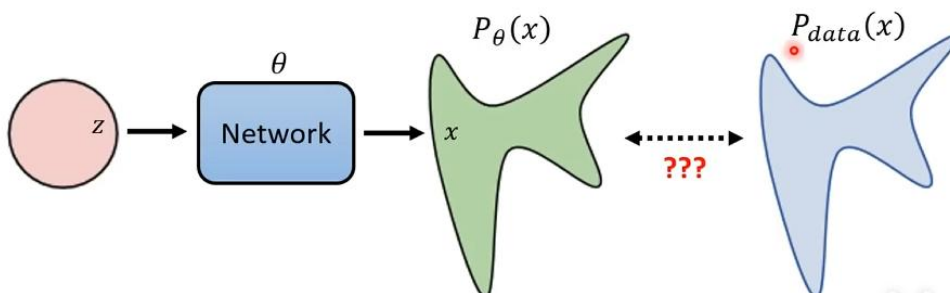
那 Sampling 过程, 也就是产生图像的过程是什么样的呢?

我们首先取一个纯粹都是噪声的图, 这个噪声符合 $(0, 1)$ 正态分布 (高斯分布), 随后循环大 T 步, 我们发现这里很奇怪的现象, 他这里又重新采样了一个高斯分布, 在这里加权放进来图片中, 也就是说, 按照我们原本的思路, 给定噪声图 x_t 与 t , 我们预测出当前的噪声, 用 x_t 减去后用这两个加权, 会得到一个相对更清晰的图片, 但是为什么这里又要加上一个噪声呢?



这是我们的第二个疑问

想要解决这个问题, 我们就必须要从长计议



回到我们刚刚讲的图像生成的任务，找到 **network**，让他能把 **z** 空间投影出来的空间与真实图像空间最一致。衡量标准是什么呢？就是我们之前所说的，找最大似然估计的方法。

给定训练数据，取样一些图片，去找一个 θ ，让他能够生成的图片 p_θ 是正确的图片的纪律最大，也就是：

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^m P_{\theta}(x^i)$$

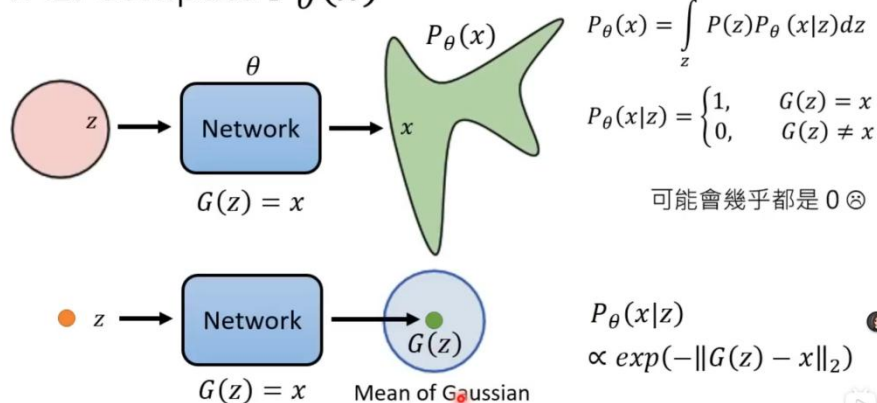
这个步骤可以写成以下形式，取 **log** 不影响，换成求和形式，变换成求期望，随后减去一个与 θ 无关的值，最后推导出的是最小散度，这就是我们的最终目标（所有的影像生成模型）

Sample $\{x^1, x^2, \dots, x^m\}$ from $P_{data}(x)$

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \prod_{i=1}^m P_{\theta}(x^i) = \arg \max_{\theta} \log \prod_{i=1}^m P_{\theta}(x^i) \\ &= \arg \max_{\theta} \sum_{i=1}^m \log P_{\theta}(x^i) \approx \arg \max_{\theta} E_{x \sim P_{data}} [\log P_{\theta}(x)] \\ &= \arg \max_{\theta} \int_x P_{data}(x) \log P_{\theta}(x) dx - \int_x P_{data}(x) \log P_{data}(x) dx \quad (\text{not related to } \theta) \\ &= \arg \max_{\theta} \int_x P_{data}(x) \log \frac{P_{\theta}(x)}{P_{data}(x)} dx = \arg \min_{\theta} KL(P_{data} || P_{\theta}) \quad \text{Difference between } P_{data} \text{ and } P_{\theta} \end{aligned}$$

那么 p_θ 如何计算？

VAE: Compute $P_{\theta}(x)$



我们可以写成上面这个西塔 z 的形式，像计算 x 的产生概率，相当于要计算给定每一个 z 产生 x 的概率的积分嘛，但是这里的给定 z 的 p 西塔 x 怎么确定呢，我们当然可以选用这个规则，但是这样计算可能会让我们结果都是 0，因为几乎不可能让生成的图片和原图完全相等，所以 VAE 给出了这样的解决办法，network 预测的这个高斯分布的均值，和对应 x 点做差取二范数，随后取复数再取对数，这样能保证大于 0 的同时，举例越小， px 越大

虽然找到了计算原则，但实际上这里的 p 西塔还是不好优化的，但它有一个下届可以最大化

VAE: Lower bound of $\log P(x)$

$$\begin{aligned}
 \log P_{\theta}(x) &= \int_z q(z|x) \log P(x) dz \quad q(z|x) \text{ can be any distribution} \\
 &= \int_z q(z|x) \log \left(\frac{P(z, x)}{P(z|x)} \right) dz = \int_z q(z|x) \log \left(\frac{P(z, x)}{q(z|x)} \cdot \frac{q(z|x)}{P(z|x)} \right) dz \\
 &= \int_z q(z|x) \log \left(\frac{P(z, x)}{q(z|x)} \right) dz + \int_z q(z|x) \log \left(\frac{q(z|x)}{P(z|x)} \right) dz \quad \geq 0 \\
 &\quad \text{KL}(q(z|x) || P(z|x)) \\
 &\geq \int_z q(z|x) \log \left(\frac{P(z, x)}{q(z|x)} \right) dz
 \end{aligned}$$

也就是说，我们对 p 西塔取 \log 以后，他一定可以写成这种这形式，这里的给定 xq 的条件概率可以是任何分布，这个等式都可以满足。具体原因比较复杂

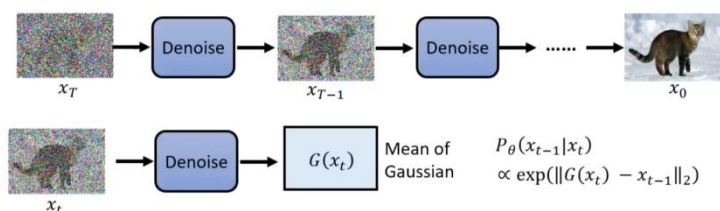
随后呢我们都学过全概率公式对吧 全概率公式就可以把这个展开成这个形式，随后天才般的构造操作出现了，上下同乘 q ，把 \log 相乘拆开，得到了这个式子

然后我们惊奇发现这个就是给定 x 的 qz 和给定 x 的 pz 的散度，散度是一定大于 0 的啊，前面这一项是可以优化的，那我们就可以优化这个下届。这一项可以写成

$$E_{q(z|x)} \left[\log \left(\frac{P(x, z)}{q(z|x)} \right) \right]$$

而对于扩散模型。

DDPM: Compute $P_{\theta}(x)$



$$P_{\theta}(x_{\theta}) = \int_{x_1: x_T} P(x_T) P_{\theta}(x_{T-1}|x_T) \dots P_{\theta}(x_{t-1}|x_t) \dots P_{\theta}(x_0|x_1) dx_1: x_T$$

p 西塔的标准和 VAE 一致，而 x 西塔，由于去噪是个马尔可夫链过程，p 西塔被写成这个形式，这样扩散模型和 vae 就被统一起来了，vae 的 q 是 decoder，而扩散模型则是 denoise

DDPM: Lower bound of $\log P(x)$

VAE Maximize $\log P_\theta(\underline{x})$ \longrightarrow Maximize $E_{q(\underline{z}|\underline{x})}[\log \left(\frac{P(\underline{x}, \underline{z})}{q(\underline{z}|\underline{x})} \right)]$

DDPM Maximize $\log P_\theta(\underline{x}_0)$ \longrightarrow Maximize $E_{q(\underline{x}_1:\underline{x}_T|\underline{x}_0)}[\log \left(\frac{P(\underline{x}_0:\underline{x}_T)}{q(\underline{x}_1:\underline{x}_T|\underline{x}_0)} \right)]$

但是

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

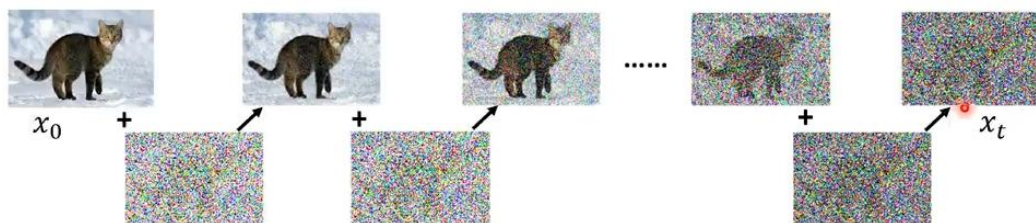
这个式子怎么被算出来呢？

我们知道 \underline{x}_{t-1} 和 \underline{x}_t 的关系是这样的， β_t 是一串预先给定的超参数，我们实际上就是在做加权

$$\begin{array}{c} \text{Image of } x_{t-1} \\ x_{t-1} \end{array} = \sqrt{1 - \beta_t} \begin{array}{c} \text{Image of } x_t \\ x_t \end{array} + \sqrt{\beta_t} \begin{array}{c} \text{Image of noise} \\ \beta_1, \beta_2, \dots, \beta_T \end{array}$$

也就是

$$q(\underline{x}_t|\underline{x}_0)$$



但实际上，这个是可以一步直接就被算出来的

$$\begin{aligned}
 x_1 &= \sqrt{1-\beta_1} x_0 + \sqrt{\beta_1} \epsilon_1 \\
 x_2 &= \sqrt{1-\beta_2} x_1 + \sqrt{\beta_2} \epsilon_2 \\
 x_2 &= \sqrt{1-\beta_2} \sqrt{1-\beta_1} x_0 + \sqrt{1-\beta_2} \sqrt{\beta_1} \epsilon_1 + \sqrt{\beta_2} \epsilon_2
 \end{aligned}$$

这里的两个噪声，其实可以简化成只去一次，并且加上这个权重

$$x_2 = \sqrt{1-\beta_2} \sqrt{1-\beta_1} x_0 + \sqrt{1-(1-\beta_2)(1-\beta_1)} \epsilon$$

$$+ \sqrt{1-(1-\beta_2)(1-\beta_1)} \epsilon$$

以此类推，整个过程可以被做成

$$x_t = \sqrt{1-\beta_1} \dots \sqrt{1-\beta_t} x_0 + \sqrt{1-(1-\beta_1) \dots (1-\beta_t)} \epsilon$$

那这里我们之前 train 算法里面的问题就被解决了

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

$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2$$
- 6: **until** converged

那我们这个也就可以化简了，不过化简的过程有点不简单啊，我们给大家拉托大的：

$$\text{Maximize } \mathbb{E}_{q(\mathbf{x}_1:\mathbf{x}_T|\mathbf{x}_0)} \left[\log \left(\frac{P(\mathbf{x}_0:\mathbf{x}_T)}{q(\mathbf{x}_1:\mathbf{x}_T|\mathbf{x}_0)} \right) \right]$$

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (47)$$

$$\begin{aligned} &= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \end{aligned} \quad (50)$$

$$= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \quad (51)$$

$$= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \quad (52)$$

$$= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \quad (53)$$

$$= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (54)$$

$$= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (55)$$

$$= \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_1:T|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (56)$$

$$= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (57)$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \quad (58)$$

化成了这个

$$\begin{aligned} & \mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T)) \\ & - \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)}[KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t))] \end{aligned}$$

第二项是不用管的，和网络的参数无关，第一项和第三项的处理方法是一样的
里面唯一不会算的东西就是给定 x_t 和 x_0 ， $q_{x_{t-1}}$ 的条件概率



我们会算的是什么呢？给定加噪声的过程，我们一步一步迭代肯定可以知道 x_{t-1} 吧，但现在这一项相当于什么？完全不给加噪声的方式，只给你加好噪声的图片和原图，让你找 x_{t-1} 的分布

通过第二托大的推导

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} \quad (71)$$

$$= \frac{\mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1-\alpha_t)\mathbf{I})\mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}x_0, (1-\bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(x_t; \sqrt{\alpha_t}x_0, (1-\alpha_t)\mathbf{I})} \quad (72)$$

$$\propto \exp \left\{ - \left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{2(1-\alpha_t)} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{2(1-\bar{\alpha}_{t-1})} - \frac{(x_t - \sqrt{\alpha_t}x_0)^2}{2(1-\alpha_t)} \right] \right\} \quad (73)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{1-\alpha_t} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0)^2}{1-\bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\alpha_t}x_0)^2}{1-\alpha_t} \right] \right\} \quad (74)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{(-2\sqrt{\alpha_t}x_t x_{t-1} + \alpha_t x_{t-1}^2)}{1-\alpha_t} + \frac{(x_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}x_{t-1}x_0)}{1-\bar{\alpha}_{t-1}} + C(x_t, x_0) \right] \right\} \quad (75)$$

$$\propto \exp \left\{ - \frac{1}{2} \left[- \frac{2\sqrt{\alpha_t}x_t x_{t-1}}{1-\alpha_t} + \frac{\alpha_t x_{t-1}^2}{1-\alpha_t} + \frac{x_{t-1}^2}{1-\bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}x_{t-1}x_0}{1-\bar{\alpha}_{t-1}} \right] \right\} \quad (76)$$

$$= \exp \left\{ - \frac{1}{2} \left[\left(\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}} \right) x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1-\bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \quad (77)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1-\bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \quad (78)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{\alpha_t - \bar{\alpha}_t + 1-\alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1-\bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \quad (79)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}x_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1-\bar{\alpha}_{t-1}} \right) x_{t-1} \right] \right\} \quad (80)$$

$$= \exp \left\{ - \frac{1}{2} \left(\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}x_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1-\bar{\alpha}_{t-1}} \right)}{\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}} x_{t-1} \right] \right\} \quad (81)$$

$$= \exp \left\{ - \frac{1}{2} \left(\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\left(\frac{\sqrt{\alpha_t}x_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}x_0}{1-\bar{\alpha}_{t-1}} \right) (1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_{t-1} \right] \right\} \quad (82)$$

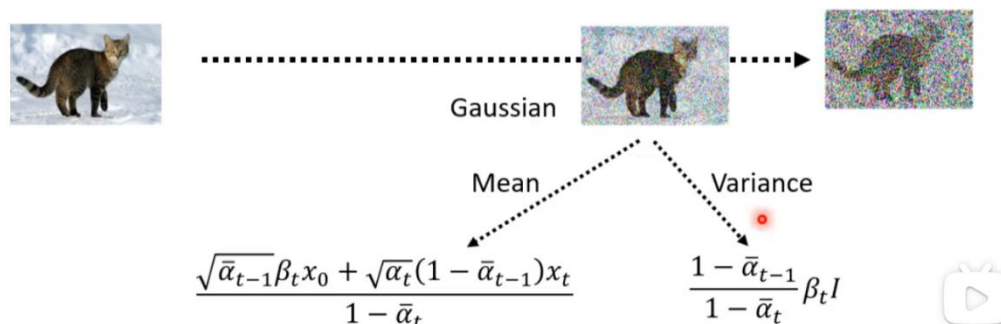
$$= \exp \left\{ - \frac{1}{2} \left(\frac{1}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[x_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)x_0}{1-\bar{\alpha}_t} x_{t-1} \right] \right\} \quad (83)$$

$$\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)x_0}{1-\bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}}_{\Sigma_q(t)} \mathbf{I}) \quad (84)$$

得到了这这也是一个高斯分布，均值和方差如下

$$E_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T E_{q(x_t|x_0)}[KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t))]$$

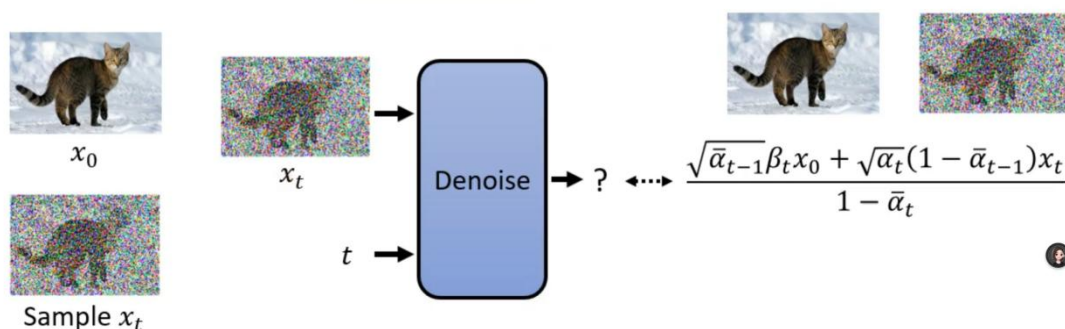


两个高斯分布的距离有具体的计算公式，也就可以极大化了

Recall that the KL Divergence between two Gaussian distributions

$$D_{KL}(\mathcal{N}(x; \mu_x, \Sigma_x) \parallel \mathcal{N}(y; \mu_y, \Sigma_y)) = \frac{1}{2} \left[\log \frac{|\Sigma_y|}{|\Sigma_x|} - d + \text{tr}(\Sigma_y^{-1} \Sigma_x) + (\mu_y - \mu_x)^T \Sigma_y^{-1} (\mu_y - \mu_x) \right]$$

最后，也就是说我们给定 x_0 , x_t 与 t ，通过去噪的过程，我们可以预测一个这个，也就是两个的加权平均（这个本质是我们刚才找到高斯分布的均值），这个继续化简，也就得到了最下面的结果



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

也就是需要预测的结果
回到

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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

取样这个算法，我们发现还有一个地方没有被解决，加这个 \mathbf{z} 干什么，本来我直接预测概率最大的高斯分布的均值就可以了啊为啥还要加扰动呢？

对于这个说法，我是没想明白的，这里给大家李彦宏老师的观点，就是说对于一个生成模型而言，如果每次都去最大概率不加扰动，就会导致生成的结果一样，但实际上我们肯定不希望如此，所以我们给每一步去噪预测，都加一个扰动，这样就可以让我生成的结果有所差异。