

# 生成扩散模型学习（一）：DDPM = 拆楼 + 建楼

朱梦

初稿于 2025-08-25，修改于 2025-08-28

## 1. 生成模型类比拆楼建楼

生成模型可以定义为将一个随机噪声  $z$  变换成一个样本数据  $x$  的过程。可以将此过程类比为“建设”，其中随机噪声  $z$  类比为砖瓦泥等建筑原材料，样本数据  $x$  类比为建筑好的高楼大厦，所以生成模型就是一支用建筑原材料建设为高楼大厦的施工队。

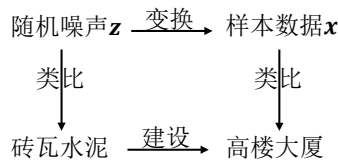


图 1 生成模型类比建楼过程

建楼过程肯定很难的，所以才用了那么多关于生成模型的研究。但是，俗话说“破坏容易建设难”，可以考虑将高楼大厦一步步地拆为砖瓦水泥的过程。设  $x_0$  为建设好的高楼大厦（即样本数据）， $x_T$  为拆好的砖瓦水泥（即随机噪声），每一步拆楼的过程定义为：

$$f_{\text{拆楼}} : x_{t-1} \rightarrow x_t \quad (1)$$

假定“拆楼”至砖瓦水泥需要  $T$  步，那么循环执行式(1)共  $T$  步，就可以得到完整的拆楼过程：

$$\text{for } i = 1 \rightarrow T : f_{\text{拆楼}} : x_{t-1} \rightarrow x_t, \text{ 即 } x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_{T-1} \rightarrow x_T \quad (2)$$

建造高楼大厦的难度在于从砖瓦水泥  $x_T$  到最终高楼大厦  $x_0$  的跨度过大，大多数人很难理解  $x_T$  是如何一步变成  $x_0$  的。但是，当有了拆楼的中间过程  $x_1, x_2, \dots, x_{T-1}$  后，那么反过来  $x_t \rightarrow x_{t-1}$  不就表示建楼的一步？如果能建模每一步建楼的过程  $f_{\text{建楼}}^{-1} : x_t \rightarrow x_{t-1}$ ，那么循环执行  $T$  步，最终不就造出高楼大厦了？

## 2. DDPM 如何拆

具体来说,去噪扩散概率模型(Denoising Diffusion Probabilistic Models,DDPM)将每一步“拆楼”的过程定义为:

$$\mathbf{x}_t = \alpha_t \mathbf{x}_{t-1} + \beta_t \epsilon_t, \text{ s.t. } \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (3)$$

其中,  $\alpha_t, \beta_t > 0$  且满足  $\alpha_t^2 + \beta_t^2 = 1$ ,  $\beta_t$  通常接近于 0。  $\beta_t \epsilon_t$  表示每一步“拆楼”过程中对“原楼体”  $\mathbf{x}_{t-1}$  的破坏程度。式(3)也可以理解为每一步“拆楼”过程中将“原楼体”  $\mathbf{x}_{t-1}$  分解为  $\alpha_t \mathbf{x}_{t-1}$  的“楼体”加上  $\beta_t \epsilon_t$  的“砖瓦水泥”

反复迭代式(3)有:

$$\begin{aligned} \mathbf{x}_t &= \alpha_t \mathbf{x}_{t-1} + \beta_t \epsilon_t \\ &= \alpha_t (\alpha_{t-1} \mathbf{x}_{t-2} + \beta_{t-1} \epsilon_{t-1}) + \beta_t \epsilon_t \\ &= \dots \\ &= (\alpha_t \dots \alpha_1) \mathbf{x}_0 + (\alpha_t \dots \alpha_2) \beta_1 \epsilon_1 + (\alpha_t \dots \alpha_3) \beta_2 \epsilon_2 + \dots + \alpha_t \beta_{t-1} \epsilon_{t-1} + \beta_t \epsilon_t \\ &= \prod_{i=1}^t \alpha_i \mathbf{x}_0 + \prod_{i=2}^t \alpha_i \beta_1 \epsilon_1 + \prod_{i=3}^t \alpha_i \beta_2 \epsilon_2 + \dots + \alpha_t \beta_{t-1} \epsilon_{t-1} + \beta_t \epsilon_t \\ &= \prod_{i=1}^t \alpha_i \mathbf{x}_0 + \sum_{j=1}^{t-1} \prod_{i=j+1}^t \alpha_i \beta_j \epsilon_j + \beta_t \epsilon_t \end{aligned} \quad (4)$$

首先,上式中  $\sum_{j=1}^{t-1} \prod_{i=j+1}^t \alpha_i \beta_j \epsilon_j + \beta_t \epsilon_t$  正好为多个独立的正太噪声之和,其均值为 0, 方差分别为  $(\alpha_t \dots \alpha_2)^2 \beta_1^2$ 、 $(\alpha_t \dots \alpha_3)^2 \beta_2^2$ 、...、 $\beta_t^2$ 。然后,根据正太分布的可叠加性,即  $\sum_{j=1}^{t-1} \prod_{i=j+1}^t \alpha_i \beta_j \epsilon_j + \beta_t \epsilon_t$  服从均值为 0、方差为  $(\alpha_t \dots \alpha_2)^2 \beta_1^2 + (\alpha_t \dots \alpha_3)^2 \beta_2^2 + \dots + \beta_t^2$  的正太分布。最后,在  $\alpha_t^2 + \beta_t^2 = 1$  恒成立之下,可以得到:

$$(\alpha_t \dots \alpha_1)^2 (\alpha_t \dots \alpha_2)^2 \beta_1^2 + (\alpha_t \dots \alpha_3)^2 \beta_2^2 + \dots + \beta_t^2 = 1 \quad (5)$$

因此,式(4)可以等价于:

$$\begin{aligned} \mathbf{x}_t &= (\alpha_t \dots \alpha_1) \mathbf{x}_0 + \sqrt{1 - (\alpha_t \dots \alpha_1)^2} \bar{\epsilon}_t \\ &\triangleq \bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \bar{\epsilon}_t, \text{ s.t. } \bar{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{aligned} \quad (6)$$

这为计算  $\mathbf{x}_t$  提供了极大的便利。且 DDPM 会选择适当的  $\alpha_t$  形式,使得  $\bar{\alpha}_T \approx 0$ , 这意味着经过  $T$  步的拆楼后,所剩的“楼体”几乎所剩无几,已经全部转换为“砖瓦水泥”。

### 3. DDPM 又如何建与损失函数

“拆楼”是  $\mathbf{x}_{t-1} \rightarrow \mathbf{x}_t$  的过程，此过程可以得到很多的数据对  $(\mathbf{x}_{t-1}, \mathbf{x}_t)$ ，那么“建楼”自然就是从这些数据对中学习一个  $\mathbf{x}_t \rightarrow \mathbf{x}_{t-1}$  的模型。首先，“拆楼”的式(3)可以改写为  $\mathbf{x}_{t-1} = (\mathbf{x}_t - \beta_t \epsilon_t) / \alpha_t$ ，这启发可以将建楼模型  $g(\mathbf{x}_t)$  设计为：

$$g(\mathbf{x}_t) = \frac{1}{\alpha_t}(\mathbf{x}_t - \beta_t \epsilon_{\theta}(\mathbf{x}_t, t)) \quad (7)$$

其中， $\theta$  表示待优化参数。容易想到的损失函数就是最小化  $\mathbf{x}_{t-1}$  和  $g(\mathbf{x}_t)$  之间的欧式距离：

$$\begin{aligned} \|\mathbf{x}_{t-1} - g(\mathbf{x}_t)\|^2 &= \left\| \frac{1}{\alpha_t}(\mathbf{x}_t - \beta_t \epsilon_t) - \frac{1}{\alpha_t}(\mathbf{x}_t - \beta_t \epsilon_{\theta}(\mathbf{x}_t, t)) \right\|^2 \\ &= \frac{\beta_t^2}{\alpha_t^2} \|\epsilon_{\theta}(\mathbf{x}_t, t) - \epsilon_t\|^2 \end{aligned} \quad (8)$$

结合式(3)、(6)有：

$$\begin{aligned} \mathbf{x}_t &= \alpha_t \mathbf{x}_{t-1} + \beta_t \epsilon_t = \alpha_t (\bar{\alpha}_{t-1} \mathbf{x}_0 + \bar{\beta}_{t-1} \bar{\epsilon}_{t-1}) + \beta_t \epsilon_t \\ &= \bar{\alpha}_t \mathbf{x}_0 + \alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t \end{aligned} \quad (9)$$

将式(9)代入式(8)，得到的损失函数形式为：

$$\frac{\beta_t^2}{\alpha_t^2} \|\epsilon_{\theta}(\bar{\alpha}_t \mathbf{x}_0 + \alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t, t) - \epsilon_t\|^2 \quad (10)$$

为什么式(9)要回退一步给出  $\mathbf{x}_t$ ？因为已经事先采样了  $\epsilon_t$ ，而  $\epsilon_t$  和  $\bar{\epsilon}_t$  并不是相互独立的，因此给定  $\epsilon_t$  的情况下，不能完全独立采样  $\bar{\epsilon}_t$ 。

### 4. 降低方差与超参设置

原则上，损失函数(10)已经可以完成 DDPM 的训练了，但它在实践中可能有方差过大的风险，从而导致收敛过慢的问题。要理解这一点并不困难，只需观察到式(10)包含了四个需要采样的随机变量：

- (1) 从所有训练样本采样一个数据样本  $\mathbf{x}_0$ 。
- (2) 从正太分布  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  中采样  $\bar{\epsilon}_{t-1}$  和  $\epsilon_t$ 。
- (3) 从  $1 \sim T$  中采样一个  $t$ 。

要采样的随机变量越多，就越难对损失函数做准确估计，即每次对损失函数进行估计的波动（方差）过大。幸运的是，可以通过积分技巧来将  $\bar{\epsilon}_{t-1}$ 、 $\epsilon_t$  合并到单个正太随机变量中，从而缓解方差过大的问题。这个积分确实有点技巧性，但也不算复杂。由于正太分布的叠加性，可得  $\alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t$  实际相当于单个随机变量  $\bar{\beta}_t \epsilon | \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 。同理， $\beta_t \bar{\epsilon}_{t-1} - \alpha_t \bar{\beta}_{t-1} \epsilon_t$  实际相当于单个随机变量  $\bar{\beta}_t \omega | \omega \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 。并且可以证明  $\mathbb{E}[\langle \epsilon, \omega \rangle] = 0$ 。接下来，用  $\epsilon$  和  $\omega$  重新表示  $\epsilon_t$

$$\epsilon_t = \frac{(\beta_t \epsilon - \alpha_t \bar{\beta}_{t-1} \omega) \bar{\beta}_t}{\beta_t^2 + \alpha_t^2 \bar{\beta}_{t-1}^2} = \frac{\beta_t \epsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t} \quad (11)$$

代入式(10)得：

$$\begin{aligned} & \mathbb{E}_{\bar{\epsilon}_{t-1}, \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \frac{\beta_t^2}{\alpha_t^2} \left\| \epsilon_{\theta}(\bar{\alpha}_t \mathbf{x}_0 + \alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t, t) - \epsilon_t \right\|^2 \right] \\ &= \frac{\beta_t^2}{\alpha_t^2} \mathbb{E}_{\epsilon, \omega \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \left\| \epsilon_{\theta}(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon, t) - \frac{\beta_t \epsilon - \alpha_t \bar{\beta}_{t-1} \omega}{\bar{\beta}_t} \right\|^2 \right] \end{aligned} \quad (12)$$

由于式(12)关于  $\omega$  只是二次的，所以可以展开将它的期望直接算出来，结果为：

$$\frac{\beta_t^2}{\alpha_t^2} \left( \frac{\beta_t^2}{\bar{\beta}_t^2} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \left\| \frac{\bar{\beta}_t}{\beta_t} \epsilon_{\theta}(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon, t) - \epsilon \right\|^2 \right] + \text{常数} \right) \quad (13)$$

省略常数和损失函数的权重，便得到 DDPM 最终所用的损失函数：

$$\left\| \frac{\bar{\beta}_t}{\beta_t} \epsilon_{\theta}(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon, t) - \epsilon \right\|^2 \quad (14)$$

DDPM 定义  $\alpha_t$  为：

$$\alpha_t = \sqrt{1 - \frac{0.02t}{T}} \quad (15)$$

在推导式(6)时有要求  $\alpha_T \approx 0$ ，大致估算：

$$\begin{aligned} \log \bar{\alpha}_T &= \sum_{t=1}^T \log \alpha_t = \frac{1}{2} \sum_{t=1}^T \log \left( 1 - \frac{0.02t}{T} \right) \\ &< \frac{1}{2} \sum_{t=1}^T \left( -\frac{0.02t}{T} \right) = -0.005(T+1) \end{aligned} \quad (16)$$

代入  $T = 1000$  大致为  $\bar{\alpha}_T = 1 \times 10^{-5}$ ，这就刚好达到约为 0 的标准。

在式(7)，显示地写出来  $t$ ，这是因为原则上不同的  $t$  处理的是不同层次的对象，所以应该使用不同的重构模型，即应该有  $T$  个不同的重构模型，于是共享所有重构模型的参数，将  $t$  作为条件传入。按照附录的说法， $t$  是转换成位置编码后，直接加到残差模块上的。

## 5. 递归生成

训练完后，可以从一个随机噪声  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  出发执行  $T$  步式(7)来进行生成。这对应于自回归解码中的 Greedy Search。如果需要进行 Random Sample，那么需要补上噪声项：

$$\mathbf{x}_{t-1} = \frac{1}{\alpha_t}(\mathbf{x}_t - \beta_t \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)) + \sigma_t \mathbf{z}, \text{ s.t. } \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (17)$$

其中， $\sigma_t = \beta_t$ ，即正向和反向的方差保持同步。从生成过程来看，生成速度为一个瓶颈，即每生成一张图像，需要反复将  $\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)$  执行  $T$  步，因此 DDPM 的一大缺陷就是采样速度慢。

## 引用源

转载于：

```
1 @online{kexuefm-9119,
2   title={生成扩散模型漫谈（一）：DDPM = 拆楼 + 建楼},
3   author={苏剑林},
4   year={2022},
5   month={Jun},
6   url={https://spaces.ac.cn/archives/9119}
7 }
```