

# 扩散模型 DDPM.

一、《一个视频看懂DDPM原理推导》——Nik-Li

人降噪过程每一步学的是一个分布(而非确切的 $X_{t-1}'$ )  $\Rightarrow$  更好随机性

但 unet 的输出最后是个噪音  $\leftarrow$  重参数化 噪音

零：理论基础。

1. 正态分布（高斯分布）， $x \sim N(\mu, \sigma^2)$

(1) 扩散过程的终点 $x_T$ 是一个标准正态分布。

(2) 多维高斯分布

$$x \sim N(\mu, \Sigma)$$

·  $x$  是图片展平后的向量

·  $\mu$ : 均值向量

·  $\Sigma$ : 协方差矩阵.  $\Rightarrow \sigma^2 I$  (每个像素噪声独立)

(3) DDPM 用到的“三大高斯性质”:

① 重参数化 (Reparameterization)

· 前向加噪核心: 让 随机采样 可导.

· 从  $N(\mu, \sigma^2)$  中采样一个  $x$ :

1° 从  $N(0, 1)$  中采样一个随机噪声  $\epsilon$ .

2° 线性变换得  $x$ :  $x = \mu + \sigma \epsilon$

· DDPM 中:

$$x_t \sim N(\sqrt{1-\beta_t} x_{t-1}, \beta_t I)$$

$$x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \cdot \epsilon \quad \begin{matrix} \leftarrow noise \\ \text{转换成预测} \end{matrix}$$

噪声  $\epsilon$ .

② 高斯分布的可加性.

# VAE 和 DDPM

(VAE, DDPM)

在生成模型论文中的符号:

- Q: 人为设定的、已知的真实分布  
· 破坏者 · 完全已知  $q(x_t | x_{t-1})$ : 在前一时刻的基础上加噪  
· Encoder

- P: 神经网络试图学习的、模拟的分布

- 重建者 · 可学习的,  $p_\theta$  ( $\theta$ 为权重)
- $p_\theta(x_{t-1} | x_t)$ : 神经网络猜上一时刻的分布
- Decoder.

• 两个独立高斯分布的随机变量相加，结果依然高斯。

$$X \sim N(\mu_1, \sigma_1^2), Y \sim N(\mu_2, \sigma_2^2) \Rightarrow X+Y \sim N(\mu_1+\mu_2, \sigma_1^2+\sigma_2^2)$$

• 在DDPM中：

让扩散中的加噪一步到位  $\Rightarrow q(x_t|x_0)$  的分布。

### ③ 贝叶斯公式与高斯乘积

$$N(x; \mu_1, \sigma_1^2) \cdot N(x; \mu_2, \sigma_2^2) \propto N(x; \mu_{\text{new}}, \sigma_{\text{new}}^2)$$

• 在DDPM中： $q(x_t|x_{t-1}) \cdot q(x_{t-1}|x_0)$ ，具体情况见下节。

## 2. 贝叶斯公式

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

(1) 先验概率 (Prior,  $P(A)$ )

• 根据以往的经验  $\rightarrow q(x_{t-1}|x_0)$

(2) 似然 (Likelihood,  $P(B|A)$ )

• 假设 A 已发生的情况下，出现 B 的概率  $\rightarrow q(x_t|x_{t-1})$

(3) 后验 (Posterior,  $P(A|B)$ ) (从果推因)

• 最想求的值。看到证据 B 后，需要修正原来的判断，得到事件 A 发生的概率。↓

逆向去噪过程：已知：全是噪声的  $x_t$  (证据 B)

想求：这图上一时刻  $x_{t-1}$  (A) 的概率

$$q(x_{t-1}|x_t)$$

(4) 贝叶斯在DDPM。

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

后验概率公式

• 假设训练集已知原因  $x_0$ :

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

• 反侧：三项全部已知且高斯

$q(x_{t-1}|x_t, x_0)$  的均值与方差可以单独解析

## 4. 联合分布 (Joint Distribution)

(1) 描述了多个随机变量同时取特定值的概率。

### (1) 离散 (Discrete)

$$P_{XY}(x,y) \triangleq P(X=x, Y=y)$$

性质：(1) 归一性

$$\sum_{x \in X} \sum_{y \in Y} P(x=x, y=y) = 1$$

- $X, Y$ : 随机变量 (不确定的函数/状态)  $P(x \leq X \leq x+dx, y \leq Y \leq y+dy) \approx f_{XY}(x,y) dx dy$
- $x, y$ : 该变量取到的具体数值.

### (2) 连续 (Continuous)

使用联合概率密度函数 (PDF):

$$f_{XY}(x,y) \rightarrow p(x,y)$$

对连续度量，单点概率  $P(x=x, y=y) = 0$   
因此我们讨论的是落在一个小区域内的概率。

归一性:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x,y) dx dy = 1.$$

## 5. 边缘分布 (Marginal Distribution)

(1) 对联合分布中“无关变量”进行求和(离散)/求积分(连续)。

### (2) 离散

$$P_X(x) = \sum_{y \in Y} P_{XY}(x,y)$$

### (3) 连续

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x,y) dy$$

↓DDPM

$$p(x) = \int p(x,y) dy.$$

## 6. 重要法则

(1) 加法法则 (边缘化)

→ 谓变量

$$P(x) = \int p(x,y) dy.$$

7. 期望 (expectation view)

· 把积分写成期望

$$P(x) = E_{Z \sim p(z)} [p(x|z)]$$

·  $p(z)$ :  $z$  服从的分布

· DDPM: 因象  $x$  的分布,本质上是无数条可能噪声路径的生成概率的平均值。

(2) 乘法

$$P(x,y) = P(y|x) \cdot P(x) = P(x|y) \cdot P(y)$$

(3) 全概率公式 → 描述隐变量模型的生成过程.

$$P(x) = \int p(x|z) \cdot p(z) dz$$

边缘概率 权重:  $z$  出现的先验概率

物理意义:  $x$  发生的概率 = 在所有可能的隐状态  $z$  下生成  $x$  的概率的加权平均.

→ 将神经网络的任务 → 学习  $q(x_{t+1} | x_t, x_0)$  的均值 → 预测噪声  $\epsilon$ .

### 3. 最大化似然 (Log likelihood) + 变分推断 (Variational Inference)

(1) 定义: 真实数据  $X = \{x_1, x_2, \dots, x_n\}$ , 模型是一个概率分布  $P_\theta(x)$ . 似然函数就是模型生成这些真实数据的联合概率:  $L(\theta) = P_\theta(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P_\theta(x_i)$   
目标是找一个  $\theta$ , 让  $L(\theta)$  最大.

(2) Why Log?

- 连乘 → 连加: 防止 Underflow
- 方便求导.

$$\Rightarrow \theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P_\theta(x_i)$$

(3) 最小化 Loss  $\Leftrightarrow$  最小化负的 Likelihood.

$$Loss = - \sum_{i=1}^n \log P_\theta(x_i)$$

(4) DDPGM 中的最大似然.

1° 由于预测数据 ~ 高斯分布  $\rightarrow P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\Rightarrow -\log P(x) = \underbrace{\frac{1}{2\sigma^2}(x-\mu)^2}_{MSE} + \underbrace{\log(\sqrt{2\pi\sigma^2})}_{\text{常数项}}$$

∴ 在高斯分布前提下: Max Log-Likelihood  $\Leftrightarrow$  MSE.

2° 由于解析的  $P_\theta(x_0)$  无法计算  $\rightarrow \int \underbrace{P_\theta(x_0, x_1, \dots, x_T)}_{\text{无数条 } x_0 \rightarrow x_1 \text{ 的路径}} dx_{1:T} = P(x_0)$  (边缘概率)

∴ 要取其下界 (Evidence Lower Bound)  $\mathbb{E}_{\text{BLBO}}$

$$\therefore ELBO \approx \sum_{\text{后验}} D_{KL}(\text{后验} || \text{模型}) \approx \sum \text{MSE}$$

高斯分布下,  $D_{KL} \propto \|\mu_1 - \mu_2\|^2 \rightarrow$  均值的均方误差.

★ 分布差异  $\Leftrightarrow$  由噪声决定

4. KL 散度 (KL Divergence)

· 衡量 2 个概率分布之间的差异.  $D_{KL}(A || B) = 0 \Leftrightarrow A \text{ 和 } B \text{ 完全一致.}$   
· 值越大, 表示差距越大.

推导:

$$X_t = \sqrt{2t} X_{t-1} + \sqrt{1-2t} \varepsilon_{t-1}$$

$$= \sqrt{2t} (\underbrace{\sqrt{2t-1} X_{t-2} + \sqrt{1-2(t-1)} \varepsilon_{t-2}}_{\text{①}}) + \sqrt{1-2t} \varepsilon_{t-1}$$

$$= \sqrt{2t} \underbrace{\sqrt{2t-1} X_{t-2}}_{\text{②}} + \underbrace{\sqrt{2t(1-2t+1)} \varepsilon_{t-2} + \sqrt{1-2t} \varepsilon_{t-1}}_{\text{noise}}$$

$$\sqrt{\text{①}^2 + \text{②}^2} = \sqrt{1-2t} \varepsilon$$

$$= \sqrt{2t} \varepsilon_{t-1} X_{t-2} + \sqrt{1-2t} \varepsilon$$

$$\Rightarrow \sqrt{2t} X_0 + \sqrt{1-2t} \varepsilon$$

$$\underbrace{(\ln \lambda)^2}_{\text{variance}} + \underbrace{(m - \bar{x})^2}_{\text{bias}} = \text{exp}(\lambda)$$

6. 白噪声  $\Leftrightarrow$  local stationary  $\Rightarrow$  均值和方差都为常数

(协方差)  $= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_{i+1} - \bar{x})$   $\rightarrow$  常数  $\Rightarrow$  方差为常数

(均值)  $\approx \frac{1}{N} \sum_{i=1}^N x_i = \bar{x} \approx 0.85$

(协方差)  $= \frac{1}{N-1} \sum_{i=1}^{N-1} (x_i - \bar{x})(x_{i+1} - \bar{x})$

均值  $\approx 0.85$  (Kolmogorov-Smirnov test)

方差  $\approx 0.01$  (Chi-squared test)

# 一、前向过程

## 1. 单步加噪:

$$x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_t \quad (\text{重参数化})$$

•  $\beta_t$ : 方差调度参数 (Variance Schedule), 很小接近于0.

•  $\varepsilon_t \sim N(0, I)$  采样的随机噪声.

## 2. 任意步加噪 ( $x_0 \rightarrow x_t$ )

• 记:  $\alpha_t = 1 - \beta_t$   $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i \quad \text{i.e. } q(x_t | x_0) \sim N(x_t; \sqrt{\bar{\alpha}_t} x_0, (I - \bar{\alpha}_t) I)$$

• 物理含义:  $\sqrt{\bar{\alpha}_t} x_0$ : 信号项. 随着  $t \uparrow$ , 原图信息↓

• 物理含义:  $\sqrt{1-\bar{\alpha}_t} \varepsilon$ : 噪声项.  $t \uparrow$ , 权重  $\rightarrow 1$ , 图像变为纯噪声.

• 当  $x_T$  中  $I$  会很大时,  $x_T \approx \varepsilon$  意味着什么?

Ans. 扩散过程最终将任意数据分布收敛到标准正态分布, 保证了我们可以用同样标准正态分布的噪声来开始生成过程.

## 二、逆向过程

$$\frac{P(x_t | x_{t-1}) \cdot \varphi(x_{t-1})}{P(x_t)} \quad \text{不好求.}$$

• 思路: 直接算  $P(x_{t-1} | x_t)$  很难, 但可以算  $q(x_{t-1} | x_t, x_0)$  的均值

i.e.  $q(x_{t-1} | x_t, x_0) \sim N(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}_t I)$

$$\tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}} \beta}{1-\bar{\alpha}_t} x_0 \quad \leftarrow \text{无法推理}$$

$\Downarrow$  任意步加噪反解  $x_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t} \varepsilon}{\sqrt{\bar{\alpha}_t}}$

$$\tilde{\mu}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \right) \varepsilon$$

$$\tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t$$

没有未知数!

•  $x_t$ : 输入给网络的当前图像 (已知)

•  $\bar{\alpha}_t, \beta_t, \beta_t$ : 人为设定超参 (已知)

•  $\varepsilon$ : 导致  $x_0$  变成  $x_t$  的那个噪声. (未知)

$\Rightarrow$  不需要预测  $x_{t-1}$ , 也不需要预测原图  $x_0$ .  
 训练网络  $E_\theta(x_t, t) \rightarrow$  预测小噪点  
 $\epsilon_\theta \rightarrow \hat{\mu}_t$  然  $x_{t-1} \rightarrow \dots \rightarrow x_0$   
 $\Downarrow \hat{\mu}_t(x_0, x_t) + \sqrt{\beta_t} \epsilon \sim N(0, 1)$

该噪声对于第7个  
 逆向步来说是分别预测  
 出来的(每步单独预测),  
 但合义时从后到前的  
 噪声.

## 2. 关于方差 $\beta_t$

- 神经网络在设定方差时直接固定.
- 物理含义:  $\sigma$ : 这步去噪有多大随机性/模糊度.  
 $\mu$ : 图像长什么样.
- IDDPM (OpenAI) 发现: 预测小学习方差, 生成的图片效果不明显, 但 log-likelihood (对数似然) 会显著提升, 且允许用更少的参数生成图片.

## 三. 常见问题.

1. 一下取多大的时候,  $x_t$  满足  $N(0, 1)$ ? Ans. 取大些.

2. Why  $\beta_t = 1 - 2t$ ? i.e.  $\mu^2 + \sigma^2 = 1$

Ans. 保证均值  $\rightarrow 0$  时, 方差  $\rightarrow 1$ .

### 3. 能否跳步?

Ans.  $P(x_{t-1}|x_t, x_0)$  是基于马尔可夫性质建立的, 所以要在逆向过程跳步的话, 需要去马尔可夫化  $\rightarrow$  DDIM.  
 (极是一种采样方式).

4. 为什么不用 UNet 去做  $x_T \xrightarrow{\text{预测}} x_{T-1}$ ?

Ans. 直接去预测  $x_{T-1}$  类似退化为风格任务. 此时每一个  $x_T$  都对应一个  $x_0$ , 即每一步都 deterministic. 取 2 个一样的  $x_T$  会得到 2 个一样的  $x_0$ .

但去做  $P(x_{t-1}|x_t)$  多样性更高.  $x_T \xrightarrow{x_0} x_0'$   
 (因为涉及采样)

#### 四、Loss Function.

$$L(\theta) = E_{t, x_0, \varepsilon} [\|\varepsilon - \varepsilon_\theta(\sqrt{a_t}x_0 + \sqrt{1-a_t}\varepsilon, t)\|^2]$$

$$LOSS = \|\varepsilon - \varepsilon_\theta(x_t, t)\|^2$$

- 上述是简化后的 Loss：将所有  $\theta$  的权重都设为 1，但简化后生成的图片质量反而更好。
- 由下界推导，我们要优化的是两个高斯分布（后验  $q$  和模型  $P$ ）之间的 KL 散度。

$$L_{\text{vib}} = E \left[ \sum_{j=1}^T \frac{\beta_t^2}{2\sigma_t^2 a_t(1-\bar{a}_t)} \|\varepsilon - \varepsilon_\theta(x_t, t)\|^2 \right]$$

- Why 简化更好？

Ans. • 加权 Loss 会更多关注  $t$  很小的步骤（纠结于 imperceptible 的细节），通常是图像几乎完全丢弃的微小噪声。  
• 不加权相当增加了  $t$  较大时的权重，迫使模型更好地学习图像的整体结构和内容。

#### 五、DDPM 中的 UNet 结构。

##### 1. DDPM 中的改动。

- 同一个 UNet 要处理  $t=1 \sim 1000$  所有情况  $\rightarrow$  需知道现在第几步。
- 方法：正弦位置编码 (Sinusoidal Positional Embeddings)。
  - 输入：标量  $t$
  - 编码：正弦/余弦  $\rightarrow$  高维向量
  - 注入：送入 UNet 的每个残差块中；

$$\text{Scale \& shift : } \text{Feature}_{\text{new}} = \text{Feature}_{\text{old}} \times (1 + \text{Scale}_t) + \text{Shift}_t$$