

扩散模型学习笔记 —— 理论

一. 预备: 1. 评价指标: IS / FID

2. ELBO + VAE + HVAG

3. VDM: $\propto \int \mathbb{E} [\log p(x)]$

二. DDPM

三. DDIM

四. SDE 统一框架

五. Analyse-DPM (expand-)

六. 概率流 ODE

七. 条件生成
| classifier guidance
| classifier-free



一、预备：

- 1) Inception Score：对于生成图片 x ，使用 Inception-V3 分类网络得到类预测 $p(y|x)$ 。

$y \in R^{1000}$

$\left\{ \begin{array}{l} p(y|x)：单一样本类别归属度 \rightarrow \downarrow \text{给出的类预测，质量高的响应样本} \\ p(y)：生成样本的多样性 \rightarrow \text{分类后统计得到，应分布均匀} \end{array} \right.$

$$\begin{aligned} \rightarrow IS(G) &= \exp \left(E_{x \sim p(x)} D_{KL}(p(y|x) || p(y)) \right) \\ &= - \underbrace{p(y|x) \log p(y)}_{\text{双端} \uparrow + p(y) \downarrow} - \underbrace{(- p(y|x) \log p(y|x))}_{\text{熵} \downarrow} \end{aligned}$$

(2) FID：为了衡量生成样本集与真实分布的距离。

将图片输入 Inception 得到分类层前的激活量 E^{2048} ， $p(y)$ 熵 \downarrow

再用高斯建模分布得到 μ 和 Σ 。

$$d^2 \left((\mu, \Sigma), (\mu_r, \Sigma_w) \right) = \|\mu - \mu_r\|^2 + \text{tr} \left(\Sigma + \Sigma_w - 2(\Sigma \Sigma_w)^{-1} \right)$$

$\underbrace{\mu \text{ 成分}}_{\text{生成分布}}, \underbrace{\Sigma \text{ 密度}}_{\text{真实分布}}$

生成模型：建模 $p(x)$ \rightarrow 可以采样 + 新样本

2. ELBO:
 (Evidence Lower
 Bound)
 $\rightarrow \log P(x)$

对于一个真实样本 x , 真实分布为 p^* . 借助隐变量 z 建模近似表示.

直接最大化似然概率: $P(x) = \int f_{P(x,z)} dz = \underbrace{\frac{P(x,z)}{P(z|x)}}_{\text{均分 } X} \text{ 未知gt编码器}$

$\log P(x) \geq E_{q_\phi(z|x)} \left[\log \frac{P(x,z)}{q_\phi(z|x)} \right] \rightarrow \text{ELBO}$
 > 可以使 $q_\phi(z|x)$ 近似 $P(z|x) \Leftrightarrow \text{最大化 ELBO}$

最小化 KL 故度 \Leftrightarrow 最大化对数似然: $E_{q_\phi(z|x)} [\ln p_0(x)] = -H(p^*) - D_{KL}(p^*(x) || p_0(x))$
 > 从真实分布中采样 $\rightarrow p_0$ 与 p^* 的相似度

$$\begin{aligned} \log P(x) &= \log P(x) \int q_\phi(z|x) dz = \int q_\phi(z|x) (\log p(x)) dz = E_{q_\phi(z|x)} [\log p(x)] \\ &= E_{q_\phi(z|x)} \left[\log \frac{P(x,z)}{q_\phi(z|x)} \right] = E_{q_\phi(z|x)} \left[\log \frac{P(x,z)}{P(z|x)} \frac{q_\phi(z|x)}{q_\phi(z|x)} \right] \\ &= E_{q_\phi(z|x)} \left[\log \frac{P(x,z)}{q_\phi(z|x)} \right] + E_{q_\phi(z|x)} \left[\log \frac{q_\phi(z|x)}{P(z|x)} \right] \\ &\geq E_{q_\phi(z|x)} \left[\log \frac{P(x,z)}{q_\phi(z|x)} \right] \quad D_{KL}(q_\phi(z|x) || P(z|x)) \geq 0 \end{aligned}$$

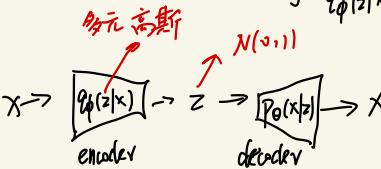
最大化 ELBO? : ① 是 $\log P(x)$ 的变分下界

② 对于给定 ϵ , ELBO + KL 故度 = 常数, 因此 $\max \text{ELBO} \rightarrow \min D_{KL}$

③ 训练后, ELBO 可以用来估计 $P(x)$ $\rightarrow q_\phi(z|x) \Rightarrow p(z|x)$

3. VAE: 变分? 求最优化 $q_\phi(z|x)$ | 自编码器? 训练时输入 = 输出.

$$E_{q_\phi(z|x)} \left[\log \frac{P(x,z)}{q_\phi(z|x)} \right] = E_{q_\phi(z|x)} \left[\log \frac{P_0(x|z)p(z)}{q_\phi(z|x)} \right] = E_{q_\phi(z|x)} [\log p_0(x|z)] - D_{KL}(q_\phi(z|x) || P(z))$$



重建向
 $z \rightarrow x$
 训练时保证
 $x \rightarrow z$ 的编码符合 $P(z)$
 防止坍缩成 Δ 函数.

① encoder: $\mu, \sigma^2 \rightarrow q_\phi(z|x) = N(z; \mu_\phi(x), \sigma^2_\phi(x)) \rightarrow$ 计算 $D_{KL}(q_\phi(z|x) || p(z))$

② 从 $q_\phi(z|x)$ 中采样 (一个样本 $z^{(l)}$) $\rightarrow E_{q_\phi}[\log p_\theta(x|z)] = \frac{1}{L} \sum_l^L \log p_\theta(x|z^{(l)})$ 解析解

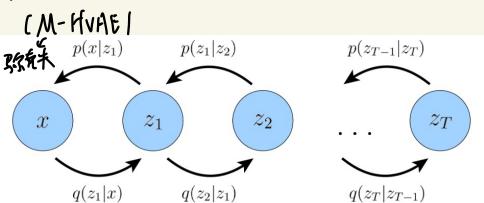
(为)使采样过程可传播度, 用重参数技巧

$$z^{(l)} = \mu_\phi(x) + \sigma_\phi(x)^2 \varepsilon, \quad \varepsilon \sim N(0, 1) \rightarrow \mu_\phi \text{ 和 } \sigma_\phi \text{ 可反传}.$$

③ 目标: 对于每个 x : $\arg\max_{\phi, \theta} \sum_{l=1}^L (\log p_\theta(x|z^{(l)}) - D_{KL}(q_\phi(z|x) || p(z)))$

[PS]: $P_\theta(x|z)$ 一般也建模为高斯, 经推导可得 $-\log p_\theta(x|z^{(l)}) = \sum_{u,v} (x^{u,v} - \mu_\theta^{u,v})^2$

4. Factor VAE. $p(x, z_t; \gamma) = p(z_T) p_\theta(x|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t)$.



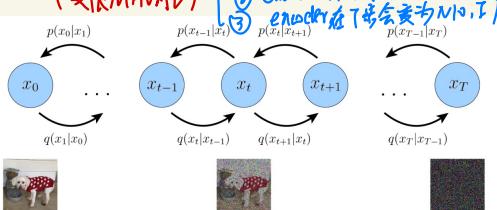
$$q_\phi(z_{1:T}|x) = q_\phi(z_1|x) \prod_{t=2}^T q_\phi(z_t|z_{t-1})$$

$$\therefore ELBO = \mathbb{E}_{q_\phi(z_{1:T}|x)} \left[\log \frac{p(x, z_{1:T})}{q_\phi(z_{1:T}|x)} \right]$$

$$= \mathbb{E}_{q_\phi} \left[\log \frac{p(z_T) p_\theta(x|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t)}{q_\phi(z_1|x) \prod_{t=2}^T q_\phi(z_t|z_{t-1})} \right]$$

5. Variational Diffusion Models: $p(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$

(类似 MHVAE) ① x 与 x_{t-1} 强度一致
② encoder 不用学, x_t 是以 x_{t-1} 为条件的高斯模型
③ encoder 在 T 时会变为 $M(0, I)$



$$q(x_t|x_{t-1}) = N(x_t; \sqrt{\kappa} x_{t-1}, (1-\kappa)t) \quad \text{加噪过程完全由高斯建模}$$

$$p(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \quad \text{保方差}$$

$$p(x_T) = N(x_T; 0, I) \quad \text{去噪训练日}$$

$$(1) \text{ ELBO} \quad \log p(x_0) = \log \int p(x_{0:T}) dx_{1:T} = \log \int \frac{p(x_{0:T}) q(x_{1:T}|x_0)}{q(x_{1:T}|x_0)} dx_{1:T} = \log \mathbb{E}_{q(x_{1:T}|x_0)} \left[\frac{p(x_{0:T})}{q(x_{1:T}|x_0)} \right]$$

$$\geq \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p(x_{0:T})}{q(x_{1:T}|x_0)} \right] = \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{\prod_{t=1}^T q(x_t|x_{t-1})} \right]$$

$$= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p(x_T) p_\theta(x_0|x_1)}{q(x_T|x_{T-1})} \right] + \mathbb{E}_{q} \left[\log \prod_{t=1}^T \frac{p_\theta(x_t|x_{t-1})}{q(x_t|x_{t-1})} \right]$$

$$\begin{aligned}
&= E_{q(x_{1:T} | x_0)} [\log p_\theta(x_0 | x_1)] + E_{q(x_{1:T} | x_0)} \left[\log \frac{p(x_T)}{q(x_{1:T} | x_0)} \right] + \sum_{t=1}^{T-1} E_{q(x_{1:t}, x_{t+1} | x_0)} \left[\log \frac{p_\theta(x_t | x_{t+1})}{q(x_{1:t}, x_{t+1} | x_0)} \right] \\
&= E_{q(x_1 | x_0)} [\log p_\theta(x_1 | x_1)] - E_{q(x_{T-1} | x_0)} [D_{KL}(q(x_T | x_{T-1}) || p(x_T))] - \sum_{t=1}^{T-1} E_{q(x_{1:t}, x_{t+1} | x_0)} [D_{KL}(q(x_t | x_{t+1}) || p_\theta(x_t | x_{t+1}))]
\end{aligned}$$

重建项，与 VAE-致
去噪匹配项
(不需优化，无参数，且由于很弱，T足够大时， $q(x_T)$ 为0)

一致项（使在前后向保持一致）
当持步长过小时有两个
随步长过大时， $T \rightarrow$ 很长，不准

另一种推导：

$$\begin{aligned}
\log p(x) &\geq E_{q(x_{1:T} | x_0)} \left[\log \frac{p(x_{1:T})}{q(x_{1:T} | x_0)} \right] = E_q \left[\log \frac{p(x_T) p_\theta(x_0 | x_1)}{q(x_1 | x_0)} \right] + \log \prod_{t=2}^T \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \\
&= E_q \left[\log \frac{p(x_T) p_\theta(x_0 | x_1)}{q(x_1 | x_0)} + \log \frac{q(x_1 | x_0)}{q(x_T | x_0)} + \log \prod_{t=2}^T \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \right] = \frac{q(x_{t-1} | x_t, x_0) q(x_t | x_0)}{q(x_{t-1} | x_0)} \\
&= E_q \left[\log p_\theta(x_0 | x_1) + E_q \left[\log \frac{p(x_T)}{p(x_T | x_0)} \right] + \sum_{t=2}^T E_q \left[\log \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \right] \right] \\
&= E_{q(x_1 | x_0)} [\log p_\theta(x_0 | x_1)] - D_{KL}(p(x_T | x_0) || p(x_T)) - \sum_{t=2}^T E_{q(x_t | x_0)} [D_{KL}(q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t))]
\end{aligned}$$

重建
去噪

去噪匹配项
 $q(x_{t-1}, x_t | x_0) = q(x_t | x_0) \cdot q(x_{t-1} | x_t, x_0)$

(当T=1时，VDM \rightarrow VAE)

① 优化：重建项如 VAE。去噪匹配项： $q(x_{t-1} | x_t, x_0) = \frac{q(x_t | x_{t-1}, x_0) q(x_{t-1} | x_0)}{q(x_t | x_0)}$

$$① q(x_t | x_{t-1}, x_0) = q(x_t | x_{t-1}) = N(x_t; \bar{x}_{t-1} x_{t-1}, \sqrt{1-\alpha_t} I)$$

$$② q(x_t | x_0) = \sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} \varepsilon_t$$

$$= \underbrace{\sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} \varepsilon_t}_{\sim N(0, I)}$$

$$= \underbrace{\sqrt{\alpha_t}_{\sum_{i=1}^T \alpha_i} x_0 + \sqrt{1-\sum_{i=1}^T \alpha_i} \varepsilon_0}_{\sim N(\bar{x}_0, \sqrt{1-\sum_{i=1}^T \alpha_i} \varepsilon_0)}$$

$$\sim N(\bar{x}_0, \sqrt{1-\sum_{i=1}^T \alpha_i} I)$$

$$\therefore q(x_{t-1} | x_t, x_0) \propto N(x_{t-1}; \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})x_t + \sqrt{1-\alpha_t}(1-\bar{\alpha}_t)x_0}{1-\bar{\alpha}_t}, \frac{(1-\bar{\alpha}_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} I)$$

$$\mu_q(x_t, x_0)$$

$$\bar{Z}_q(t)$$

$\therefore q(x_{t-1} | x_t, x_0)$ 也是正态分布，因此 $P_\theta(x_{t-1} | x_t)$ 直接建模，得到 $\mu_\theta(t)$, $\sigma_\theta^2(t)$

由未观测到 x_{t-1} 时参数确定，
经建模， $\sigma_\theta^2(t) = \bar{Z}_q(t)$

$$\underset{\theta}{\operatorname{argmin}} D_{KL}(q(x_{t-1} | x_t, x_0) \| P_\theta(x_{t-1} | x_t))$$

$$= \underset{\theta}{\operatorname{argmin}} D_{KL}\left(N(x_{t-1}; \mu_\theta, \bar{Z}_q(t)) \| N(x_{t-1}; \mu_\theta, Z_q(t))\right)$$

解解解

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma_\theta^2(t)} \left[\|\mu_\theta - \mu_q\|_2^2 \right]_{\mu_q(x_t, x_0)}$$

$$\text{将 } \mu_\theta \text{ 与 } \mu_q \text{ 表示: } \mu_\theta(x_t, t) = \frac{\bar{x}_{t-1}(1-\bar{\alpha}_t) + \bar{\alpha}_{t-1}(1-\bar{\alpha}_t)\hat{x}_\theta(x_t, t)}{1-\bar{\alpha}_t}$$

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma_\theta^2(t)} \frac{\bar{\alpha}_{t-1}(1-\bar{\alpha}_t)^2}{(1-\bar{\alpha}_t)^2} \left[\|\hat{x}_\theta(x_t, t) - x_0\|_2^2 \right] \text{ 相当于, 训一个 MN. 可以以帧为输入预测原始图片。}$$

$$\xrightarrow{\text{沿时间求和}} \underset{\theta}{\operatorname{argmin}} E_{t \sim U\{1, T\}} \left[\bar{E}_{q(x_{t-1} | x_0)} \left[D_{KL}(q(x_{t-1} | x_t, x_0) \| P_\theta(x_{t-1} | x_t)) \right] \right]$$

→ 可以沿时间步随机采样来优化。

(3) 噪声参数的优化：用 MN 建模 $\bar{\alpha}_t$ 是低效的，因为推理由需要计算 $\bar{\alpha}_t$

$$\text{目标: } \frac{1}{2\sigma_\theta^2(t)} \frac{\bar{\alpha}_{t-1}(1-\bar{\alpha}_t)^2}{(1-\bar{\alpha}_t)^2} \left[\|\hat{x}_\theta(x_t, t) - x_0\|_2^2 \right] = \frac{1}{2} \left(\frac{\bar{\alpha}_{t-1}}{1-\bar{\alpha}_{t-1}} - \frac{\bar{\alpha}_t}{1-\bar{\alpha}_t} \right) \left[\|\hat{x}_\theta(x_t, t) - x_0\|_2^2 \right]$$

$$= \frac{1}{2} (\text{SNR}(t-1) - \text{SNR}(t)) \left[\|\hat{x}_\theta(x_t, t) - x_0\|_2^2 \right] \quad \text{SNR} := \frac{\mu^2}{\sigma^2}, q(x_t | x_0) \\ = (x_t, \sqrt{1-\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) I)$$

→ 直接用 MN 建模 SNR.

表示原信号(μ)与噪声量的比率

$$\text{SNR}(t) = e^{-w_g(t)} \rightarrow \text{递增} \quad (\text{由于随 t 推进, SNR 变小, 直到 } 0, N(0, I))$$

$$\rightarrow \bar{\alpha}_t = \text{sigmoid}(-w_g(t)), 1-\bar{\alpha}_t = \text{sigmoid}(w_g(t))$$

$$\frac{1}{\bar{\alpha}-1} - \frac{1-\bar{\alpha}}{1-\bar{\alpha}} \quad \frac{1}{1-\bar{\alpha}} - 1 + \frac{1-\bar{\alpha}}{1-\bar{\alpha}} \quad \frac{1}{\bar{\alpha}(1-\bar{\alpha})^2}$$

(4) 其它解释：① VDM \Leftrightarrow 从 x_t 中预测原始图片 x_0

② \Leftrightarrow 预测 x_0 ，DDPM 但是如此，效果更好！

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \varepsilon \rightarrow x_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t} \varepsilon}{\sqrt{\bar{\alpha}_t}}$$

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

预测 $x_t \rightarrow x_0 \rightarrow \varepsilon$

$$= \frac{1}{\sqrt{\alpha t}} X_t - \frac{1-\bar{\alpha}t}{\sqrt{1-\bar{\alpha}t} \sqrt{\alpha t}} \underbrace{\Sigma_0}_{\text{S}_0}$$

$$\text{写成相册形式 } \mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha t}} X_t - \frac{1-\bar{\alpha}t}{\sqrt{1-\bar{\alpha}t} \sqrt{\alpha t}} \underbrace{\Sigma_0}_{\text{S}_0}(x_t, t)$$

$$\therefore \text{目标} \Rightarrow \frac{1}{2\sigma_q^2(t)} \frac{(1-\bar{\alpha}t)^2}{(1-\bar{\alpha}t)\alpha t} \left[\|\Sigma_0 - \hat{\Sigma}_0(x_t, t)\|_F^2 \right]$$

③ \hookrightarrow 预测得分函数 (在数据空间的梯度), 为噪音水平, $\nabla \log P(x_t)$

从 Tweedie's Formula 角度: 对一个指数分布族的均值, 可以采样样本, 再经由平均值 + 涉及分数的校正项来估计
(若采样即为本身)

$$q(x_t | x_0) = N(x_t; \sqrt{\alpha t} x_0, (1-\bar{\alpha}t) I)$$

$$\rightarrow E[\mu_{\theta}(x_t | x_0)] = x_t + (1-\bar{\alpha}t) \nabla_{x_t} \log P(x_t)$$

作为分子, 会将往 P 增大的方向修正

$$\therefore x_0 = \frac{x_t + (1-\bar{\alpha}t) \nabla \log P(x_t)}{\sqrt{\alpha t}} = \frac{x_t - \sqrt{1-\bar{\alpha}t} \Sigma_0}{\sqrt{\alpha t}} \Rightarrow \nabla \log P(x_t) = \frac{-1}{\sqrt{\alpha t}}$$

$$\text{代入 } \mu_{\theta}(x_t, x_0) = \frac{1}{\sqrt{\alpha t}} X_t + \frac{1-\bar{\alpha}t}{\sqrt{\alpha t}} \nabla \log P(x_t)$$

使 x_t 的 $\log P$ 增大的方向,
即“噪音”的反方向.

$$\text{将 } \mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha t}} X_t + \frac{1-\bar{\alpha}t}{\sqrt{\alpha t}} \underbrace{\Sigma_0}_{\text{S}_0}(x_t, t) \text{ 预测分子}$$

$$\text{目标: } \underset{\theta}{\arg \min} \frac{1}{2\sigma_q^2(t)} \frac{(1-\bar{\alpha}t)^2}{\alpha t} \left[\|\Sigma_0(x_t, t) - \nabla \log P(x_t)\|_F^2 \right]$$

6. Score-based Generative Models: 建模 score function, 并使用梯度下降特卡罗法 (如翻之万动力学) 生成样本

$$\text{任何概率分布均可以写成 } p_{\theta}(x) = \frac{1}{Z_{\theta}} e^{-f_{\theta}(x)} \rightarrow \text{能量函数}$$

由于 f_{θ} 复杂 \rightarrow 无法极大似然学习 $p_{\theta}(x)$ \rightarrow 用 $\nabla_x f_{\theta}(x)$

$$\rightarrow \nabla_x \log p_{\theta}(x) = \nabla_x \log \frac{1}{Z_{\theta}} e^{-f_{\theta}(x)}$$

$$= \nabla_x \log \frac{1}{Z_{\theta}} + \nabla_x \log e^{-f_{\theta}(x)} = -\nabla_x f_{\theta}(x) \propto S_{\theta}(x)$$

$$\rightarrow \text{最小化 } E_{p(x)} \left[\|\Sigma_0(x) - \nabla \log p_{\theta}(x)\|_F^2 \right] \rightarrow \text{得分匹配 loss, 通过对 } \nabla \log p_{\theta}(x) \text{ 的需求而最小化 KL 散度}$$

得分函数 $\nabla \log p(x)$ 的含义: 给定 x 的 \log 似然的梯度 \rightarrow 给出使似然增大的方向

从而, 通过学习 S_{θ} , 可以从样本空间任意点开始, 迭代地根据 score 更新位置来产生样本.

$$x_{t+1} \leftarrow x_t + \nabla \log p(x_t) + \sqrt{\alpha t} S_{\theta}$$

翻之万动力学

$$\text{二. DDPM: } L = E_{q_t} [D_{KL}(q(x_t | x_0) || p(x_t)) + \sum_{t>1} D_{KL}(q(x_{t+1} | x_t, x_0) || p_\theta(x_{t+1} | x_t)) - \underbrace{\log p_\theta(x_0 | x_1)}_{L_0}]$$

L_t

L_{t+1}

简陋版

L_0

1. 前向: α_t / β_t 固定, 无可学习参数, L_t 忽略.

$$= E_{x_t, x_0, \Sigma} [\| \Sigma - \Sigma_0 (\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \Sigma, t) \|_F^2]$$

2. 反向: $p_\theta(x_{t+1} | x_t) = N(x_{t+1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$

将 $\mu_\theta(x_t, t)$ 参数化为 $\frac{1}{\sqrt{\alpha_t}} (x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \Sigma_0(x_t, t))$ set to b^2 [天黑训练] 七小时前的权重下降(应该得把优化目标量化一下), 从而使模型得注意力放到长时.

$$\Rightarrow L_t = \frac{1}{\alpha_t} \frac{\beta_t^2}{1-\alpha_t} \frac{1}{\alpha_t} \frac{\partial}{\partial t} (\sqrt{1-\alpha_t}) \| \Sigma - \Sigma_0(x_t, t) \|_F^2$$

train: ∇_{x_0} , 采样 $t \sim U[1, T]$, $\Sigma \sim N(0, I)$

$$\rightarrow x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \Sigma$$

$$\rightarrow L = \| \Sigma - \Sigma_0(x_t, t) \|_F^2$$

sample (infer): $\# x_T \sim N(0, I)$.

for $t = T, \dots, 1$:

$$x_{t+1} = \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \Sigma_0(x_t, t)) + b_t \Sigma \rightarrow z \sim N(0, I), t > 1$$

前向过程不再引入 x_0 , 而均与 x_{t+1}, x_0 独立 $\mu_\theta(x_t, t)$

$$\text{三. DDIM: } q_b(x_{t+1} | x_0) := q_b(x_{t+1} | x_0) \prod_{t=2}^T q_b(x_{t+1} | x_t, x_0) \quad \text{Ayg: DDPM: } q(x_t | x_{t+1}) \rightarrow q(x_t | x_0)$$

①从更高角度, 将 DDPM 的训练过程由逐帧转换为推倒非逐帧 $q(x_{t+1} | x_t, x_0) = \frac{q(x_t | x_{t+1}, x_0) q(x_{t+1} | x_0)}{q(x_t | x_0)}$ (逐帧)

(并未改变目标, 仅是等价解/解释) \downarrow 提供只依赖 $q(x_t | x_0)$ \rightarrow 省去 $q(x_t | x_{t+1})$

②对比, 提出一种加速采样方法. $q_b(x_{t+1} | x_0) = \int q_b(x_{t+1} | x_t, x_0) q(x_t | x_0) dx_t$ 从假设 $q(x_t | x_{t+1}) \rightarrow q(x_t | x_0)$

$$N(\sqrt{\alpha_{t+1}} x_0, (1-\alpha_{t+1}) I) \downarrow \rightarrow N(\sqrt{\alpha_t} x_0, (1-\alpha_t) I)$$

待定系数可求解: $x_{t+1} = k_t (\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \Sigma) + \lambda_t x_0 + b_t \Sigma$

$$N(x_{t+1}; k_t \sqrt{\alpha_t} x_0 + \lambda_t x_0, b_t^2 I) = (k_t \sqrt{\alpha_t} + \lambda_t) x_0 + (k_t \sqrt{1-\alpha_t} \Sigma + b_t \Sigma)$$

得到的 x_{t+1} 相当于采样到 x_t . 再根据 x_t 和 x_{t+1} 的分布, 此时再只关注 x_t 的分布, 即 $\therefore k_t \sqrt{\alpha_t} + \lambda_t = \sqrt{\alpha_{t+1}}, \sqrt{k_t^2 + (1-\alpha_t)} = \sqrt{1-\alpha_{t+1}}$ 相当于通过积分削去 x_t .

$$\therefore q_b(x_{t+1} | x_t, x_0) = N(\sqrt{\alpha_{t+1}} x_0 + \sqrt{1-\alpha_{t+1}-b_t^2} \cdot \frac{x_t - \sqrt{\alpha_t} x_0}{\sqrt{1-\alpha_t}}, b_t^2 I) \quad \begin{cases} k_t, \lambda_t, b_t \\ \text{只有两个方程} \end{cases} \rightarrow \text{解方程}$$

后与 DDPM 类似. 代入 $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \Sigma$

$$\therefore \mu_q(x_{t+1} | x_t, x_0) = \frac{x_t - \sqrt{\alpha_t} x_0}{\sqrt{1-\alpha_t}} + \sqrt{1-\alpha_{t+1}-b_t^2} \Sigma$$

$$= \frac{x_t}{\sqrt{\alpha_t}} + (\sqrt{1 - \frac{\alpha_t}{\alpha_{t-1}}} - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_{t-1}}}) \underbrace{\varepsilon_0(x_t, t)}_{\text{随机噪声}}$$

$$\therefore \mu_\theta(x_{t+1} | x_t) = \frac{1}{\sqrt{\alpha_t}} x_t + \frac{1}{\sqrt{\alpha_t}} \left(\sqrt{\alpha_t (1 - \frac{\alpha_t}{\alpha_{t-1}} - \frac{\alpha_t}{\alpha_{t-1}})} - \sqrt{1 - \alpha_t} \right) \underbrace{\varepsilon_0(x_t, t)}_{\text{随机噪声}}$$

α_t 决定了前向过程的随机程度，当 $\alpha_t \rightarrow 0$ 时，只要状态 x_t 和 t 确定， x_{t+1} 就确定。

$$p_\theta(x_{t+1} | x_t) = \begin{cases} N \left(\frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1 - \alpha_t} \varepsilon_0), \sigma^2 \right), & t=1 \\ q_\theta(x_{t+1} | x_t, x_0 = \frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1 - \alpha_t} \varepsilon_0)) & t>1 \end{cases}$$

当 $\alpha_t = \frac{\sqrt{1 - \alpha_{t-1}}}{1 - \alpha_{t-1}} \sqrt{1 - \alpha_t}$ 时，等价于马尔科夫的 DDPM.

当 $\alpha_t = 0$ ，DDIM（一个以 DDPM 训练目标的隐含概率模型）
implicit

加速采样：由于前向过程可以视为非马尔科夫过程。

对于一个训练的采样时，可以将其作为 $\dim T$ 步训练好的模型。 $T = [T_1, \dots, T_{\dim(T)}]$ 为 $[1, T]$ 的 T 步 DDPM 模型，
→ 生成也只需 $\dim T$ 步，任意子序列。

$$p(x_{T_i-1} | x_{T_i}) = N \left(x_{T_i-1} ; \frac{\bar{\alpha}_{T_i-1}}{\bar{\alpha}_{T_i}} \left(x_{T_i} - \left(\frac{\bar{\beta}_{T_i}}{\bar{\alpha}_{T_i-1}} \sqrt{\bar{\beta}_{T_i}^2 - \bar{\beta}_{T_i}^2} \right) \varepsilon_\theta(x_{T_i}, T_i) \right) \right)$$

文中给出了两种方法：线性采样 / 二次采样

II

DDIM 没改变 DDPM 的训练，而仅改变了采样方法

drift 系数 diffusion 系数

I: 后进

II: $\alpha_t = 0$ 时步隐含性 DDIM

四、Score-based: DDPM: 离散的 T 步 → SDE: 时间上的连续微分方程

1. 离散 SDE 框架：前向进程： $dx = f_t(x) dt + g_t dw$ 布朗运动
 (有确定的 f_t 和 g_t) 可视作离散过程： $x_{t+\Delta t} - x_t = f_t(x_t) \Delta t + g_t \sqrt{\Delta t} \xi$, 当 $\Delta t \rightarrow 0$ 时的极限
 .. 越小的 Δt , 代表对 SDE 越好的近似 类似于 $f_t(x) \rightarrow dx/dt$

DDPM 中不同的 T , 可视为对 SDE 不同离散化程度的体现。

→ 引入 SDE 可以“在分析时借助连续性 SDE”，而在实践时适当的离散化。

逆向进程： $dx = [f_t(x) - g_t^2 \nabla_x \log p_t(x)] dt + g_t dw$ (推导省去)
 未知，“score”

→ 离散: $x_t - x_{t+1} = -[f_{\text{prior}}(x_{t+1}) - g_t^2 \nabla_{x_{t+1}} \log p(x_{t+1})] \Delta t + g_{t+1} \Delta t \xi$

$$p(x_t) = \int p(x_t | x_0) \tilde{p}(x_0) dx_0 = E_{x_0} [p(x_t | x_0)]$$

$$\rightarrow \nabla_{x_t} \log p(x_t) = \frac{\nabla_{x_t} p(x_t)}{p(x_t)} = \frac{E_{x_0} [\nabla_{x_t} p(x_t | x_0)]}{E_{x_0} [p(x_t | x_0)]} = \frac{E_{x_0} [\nabla_{x_t} p(x_t | x_0) \nabla_{x_t} \log p(x_t | x_0)]}{E_{x_0} [p(x_t | x_0)]}$$

To 使模型强, 则 $S_\theta(x_t, t) \approx \nabla_{x_t} \log p(x_t)$

可视为以 $p(x_t | x_0) p(x_0)$ 为权重对

$$\min E_{x_0} [p(x_t | x_0) \| S_\theta(x_t, t) - \nabla_{x_t} \log p(x_t | x_0) \|^2]$$

对 x_t 求期望以对每个 x_t 最小化
可忽略仅加权 loss
 $E_{x_0} [p(x_t | x_0)] \rightarrow$ 所谓, 若 $f_t(x)$ 为仿射变换, 则 $p(x_t | x_0)$ 也是高斯

$$\rightarrow \min E_{x_0, x_t} [p(x_t | x_0) \tilde{p}(x_0) \int \| S_\theta(x_t, t) - \nabla_{x_t} \log p(x_t | x_0) \|^2] \text{ 得分匹配的损失函数}$$

2. 解释DDPM: $p(x_t | x_{t+1}) = N(x_t; \sqrt{\alpha_t} x_{t+1}, (1-\alpha_t) I) \rightarrow x_t = \sqrt{\alpha_t} x_{t+1} + (1-\alpha_t) \xi$

$$\Rightarrow SDE: dx = -\frac{1}{2} \beta(t) x dt + \sqrt{\beta(t)} dw$$

[推导: 定义一个新的序列: $\{\bar{\beta}_i\}_{i=1}^N$ ($\beta_i := 1 - \alpha_i$)

$$\rightarrow x_t = \sqrt{\frac{\bar{\beta}_1}{N}} x_{t+1} + \sqrt{\frac{\bar{\beta}_1}{N}} \xi, \text{ 当 } N \rightarrow \infty \text{ 时, } \{\bar{\beta}_i\}_{i=1}^N \xrightarrow{t \in [0, 1]} \beta(t), \bar{\beta}\left(\frac{i}{N}\right) = \bar{\beta}_i, x\left(\frac{i}{N}\right) = x_i;$$

$$\stackrel{\Delta t = \frac{1}{N}}{\rightarrow} x(t+\Delta t) = \sqrt{1 - \beta(t+\Delta t)} \Delta t x(t) + \sqrt{\beta(t+\Delta t)} \Delta t \xi$$

注: $q(x_t | x_0) = N(x_t; \alpha x_0, \beta t^2 I)$ $\approx (1 - \frac{1}{2} \beta(t+\Delta t) \Delta t) x(t) + \sqrt{\beta(t+\Delta t)} \sqrt{\Delta t} \xi$

$$\Leftrightarrow f_t = \frac{d \log q_t}{dt} \approx \left(-\frac{1}{2} \beta(t) \Delta t \right) x(t) + \sqrt{\beta(t)} \sqrt{\Delta t} \xi$$

$$g_t^2 = \frac{d \sigma_t^2}{dt} = \frac{2 d \log q_t}{dt} \approx \frac{f_t}{g_t} = x(t) - \underbrace{\frac{1}{2} \beta(t) x(t) \Delta t}_{f_t} + \underbrace{\frac{\sqrt{\beta(t)}}{\sqrt{\Delta t}} \sqrt{\Delta t} \xi}_{g_t}$$

五. Analytic DPM: 在DDIM的基础上推出了方差的解析解.

$$\text{DDIM: } p(x_{t+1} | x_t, x_0) = N\left(x_{t+1}; \frac{\sqrt{\bar{\beta}_t} - \bar{\beta}_t^2}{\bar{\beta}_t} x_t + \underbrace{\left(\frac{1}{\bar{\beta}_t} - \frac{\bar{\beta}_t \sqrt{\bar{\beta}_t} - \bar{\beta}_t^2}{\bar{\beta}_t}\right)}_{Y_t} x_0, \bar{\beta}_t^2 I\right)$$

$$\rightarrow p(x_{t+1} | x_t) \approx p(x_{t+1} | x_t, x_0 = \bar{p}(x_t))$$

$$p(x_{t+1} | x_t) = \int p(x_{t+1} | x_t, x_0) p(x_0 | x_t) dx_0$$

→ 用 $N(x_0; \bar{m}(x_t), \bar{\sigma}^2 I)$ 逼近 $p(x_0 | x_t)$, 则

$$p(x_{t+1} | x_t) = N\left(x_{t+1}; \underbrace{\frac{\sqrt{\bar{\beta}_t} - \bar{\beta}_t^2}{\bar{\beta}_t} x_t + \bar{x}_t \bar{p}(x_t)}_{\text{与过去一致}}, \underbrace{(\bar{\sigma}^2 + \bar{\sigma}^2 \bar{y}_t^2)}_{\neq \bar{\sigma}^2, \bar{y}_t^2 \text{ 即最优化的修正项.}} I\right)$$

优化：待预测 $\bar{\mu}(x_t)$, \bar{b}_t

均值优化不变，将 $\bar{\mu}(x_t)$ 参数化后得到 $\|\xi - \xi_\theta(x_t, t)\|^2$

$$\bar{\mu}(x_t) = \frac{1}{\alpha_t} (x_t - \bar{b}_t \xi_\theta(x_t, t))$$

$$\begin{aligned}
 \text{证: } \mathbb{E}(x_t) &= \mathbb{E}_{x_0 \sim p(x_0|x_t)} [(x_0 - \bar{\mu}(x_t))(x_0 - \bar{\mu}(x_t))^T] \\
 &= \mathbb{E}_{x_0} \left[\left((x_0 - \frac{x_t}{\alpha_t}) + \frac{\bar{b}_t}{\alpha_t} \xi_\theta(x_t, t) \right) \left((x_0 - \frac{x_t}{\alpha_t}) + \frac{\bar{b}_t}{\alpha_t} \xi_\theta(x_t, t) \right)^T \right] \\
 &= \mathbb{E}_{x_0} \left[\left(x_0 - \frac{x_t}{\alpha_t} \right) \left(x_0 - \frac{x_t}{\alpha_t} \right)^T \right] - \frac{\bar{b}_t^2}{\alpha_t^2} \xi_\theta(x_t, t) \xi_\theta(x_t, t)^T \\
 &= \frac{1}{\alpha_t^2} \mathbb{E}_{x_0 \sim p(x_0|x_t)} x_t - \bar{b}_t x_0^T - \frac{\bar{b}_t^2}{\alpha_t^2} \xi_\theta(x_t, t) \xi_\theta(x_t, t)^T
 \end{aligned}$$

我们希望 b_t 与 x_t 无关 \rightarrow 对 $x_t \sim p(x_t)$ 求均值。

$$\begin{aligned}
 &\mathbb{E}_{x_t \sim p(x_t)} \mathbb{E}_{x_0 \sim p(x_0|x_t)} [(x_t - \bar{b}_t x_0)(x_t - \bar{b}_t x_0)^T] \\
 &= \mathbb{E}_{x_0 \sim p(x_0)} \mathbb{E}_{x_t \sim p(x_t|x_0)} [(x_t - \bar{b}_t x_0)(x_t - \bar{b}_t x_0)^T] \quad \text{关键} \\
 &= \mathbb{E}_{x_0 \sim p(x_0)} [\bar{b}_t^2 I] = \bar{b}_t^2 I \quad p(x_t|x_0) \text{ 的 } \mu \quad p(x_t|x_0) \text{ 的 } \Sigma \\
 &\therefore \mathbb{E}_t = \mathbb{E}_{x_t \sim p(x_t)} [\mathbb{E}(x_t)] = \frac{\bar{b}_t^2}{\alpha_t^2} (I - \mathbb{E}_{x_t \sim p(x_t)} [\xi_\theta(x_t, t) \xi_\theta(x_t, t)^T]) \\
 &\text{能为常量} \quad \bar{b}_t^2 = \frac{\bar{b}_t^2}{\alpha_t^2} \left(1 - \frac{1}{d} \mathbb{E}_{x_t \sim p(x_t)} [\|\xi_\theta(x_t, t)\|^2] \right) \leq \frac{\bar{b}_0^2}{\alpha_t^2} \quad (\text{原文是用得前面数据的})
 \end{aligned}$$

累积均值 $\frac{1}{M} \sum_{m=1}^M \|\xi_\theta(x_{n,m}, t)\|^2$
 估计
 ↓
 时间步

Extended-Analytic-DPM: 之前假设 $p(x_0|x_t)$ 的协方差为 $b_t^2 I$ ，即对角阵且元素相等

四、extend: 概率流 ODE: DDM 的连续形式 SDE 离散时步长不能太大

1. F-P 方程: 对于分布 $p(x) = \int \Delta(x-y) p(y) dy = \mathbb{E}_y [\Delta(x-y)]$ 转为期望

$$\rightarrow p(x)f(x) = \int \Delta(x-y) p(y) f(y) dy = \mathbb{E}_y [\Delta(x-y) f(y)]$$

$$\rightarrow \nabla_x [p(x)f(x)] = \mathbb{E}_y [\nabla_x \Delta(x-y) f(y)] = \mathbb{E}_y [f(y) \nabla_x \Delta(x-y)] \quad \textcircled{1}$$

2. SDE: $dx = f_t(x) dt + g_t dw \rightarrow x_{t+\Delta t} = x_t + f_t(x_t) \Delta t + g_t \sqrt{\Delta t} \xi, \xi \sim N(0, 1)$

$$\therefore \Delta(x - x_{t+\Delta t}) = \Delta(x - x_t - f_t(x_t) \Delta t - g_t \sqrt{\Delta t} \xi)$$

$$\approx \Delta(x - x_t) - (f_t(x_t) \Delta t + g_t \sqrt{\Delta t} \xi) \nabla_x \Delta(x - x_t) + \frac{1}{2} (g_t \sqrt{\Delta t} \xi)^2 \nabla_x^2 \Delta(x - x_t) + O(\Delta t)$$

左加求期望可得：

$$P_{t+\Delta t}(x) = E_{x_{t+\Delta t}}(\Delta(x - x_t + \Delta t)) \approx E_{x_t, \xi} [\Delta(x - x_t) - (f_t(x_t) \Delta t + g_t \sqrt{\Sigma} \xi)] \nabla_x \Delta(x - x_t) + \frac{1}{2} (g_t^T \Sigma g_t)$$

$E(\xi) = 0$ $B(\xi^2) = I$

$$= E_{x_t} [\Delta(x - x_t) - f_t(x_t) \Delta t \cdot \nabla_x \Delta(x - x_t) + \frac{1}{2} g_t^2 \Delta t^2 \nabla_x^2 \Delta(x - x_t)]$$

由①

$$= P_t(x) - \nabla_x [f_t(x) \Delta t] P_t(x) + \frac{1}{2} g_t^2 \Delta t \nabla_x^2 P_t(x)$$

原以 Δt 并 \rightarrow

$$\Rightarrow \frac{\partial P_t(x)}{\partial t} = -\nabla_x [f_t(x) P_t(x)] + \frac{1}{2} g_t^2 \nabla_x^2 P_t(x)$$

原 SDE 的 F-P 方程，描述边际分布的微分方程。

2. 等价变换：对 $\# \Delta t^2 \leq g_t^2$ ，
 $\frac{\partial P_t(x)}{\partial t} = -\nabla_x [f_t(x) P_t(x) - \frac{1}{2} (g_t^2 - b_t^2) \nabla_x P_t(x)] + \frac{1}{2} b_t^2 \nabla_x^2 P_t(x)$

$$= -\nabla_x [(f_t(x) - \frac{1}{2} (g_t^2 - b_t^2) \nabla_x \log P_t(x)) P_t(x)] + \frac{1}{2} b_t^2 \nabla_x^2 P_t(x)$$

相当于 $g_t \rightarrow b_t$ ，而 $f_t \rightarrow f_t(x) - \frac{1}{2} (g_t^2 - b_t^2) \nabla_x \log P_t(x)$

→ 将 F-P 方程对应回 SDE：

$$dx = (f_t(x) - \frac{1}{2} (g_t^2 - b_t^2) \nabla_x \log P_t(x)) dt + b_t dw$$

完全等价于原 SDE

→ 存在 不同方差的前向过程，产生的边际分布是一样的。→ DDM 可统一起来

$$\rightarrow 反向 SDE: dx = (f_t(x) - \frac{1}{2} (g_t^2 + b_t^2) \nabla_x \log P_t(x)) dt + b_t dw$$

3. 神经 ODE：当取 $b_t = 0$ 时，SDE 退化为 ODE (常微分)

$$dx = (f_t(x) - \frac{1}{2} g_t^2 \nabla_x \log P_t(x)) dt$$

概率流 ODE
用 $S_\theta(x, t)$ 近似 → 对应一个“神经 ODE”

此时，反向与前向 ODE 完全相同，可逆变换。

且前向与反向均为确定性变换。

没有随机性，从而可以在高维时
取更大步长

六、条件生成：

$$\nabla (\log p(x_t | y)) = \nabla \log \left(\frac{p(x_t) p(y|x_t)}{p(y)} \right) = \nabla \log p(x_t) + \nabla \log p(y|x_t)$$

1. Classifier-Guidance: (只需求 classifier)

$$(1) p(x_{t+1} | x_t) \xrightarrow{\text{梯度}} p(x_{t+1} | x_t, y) = \frac{p(x_{t+1} | x_t, y) p(y|x_t, x_{t+1})}{p(y|x_t)}$$

$\frac{p(x_{t+1}|y)}{p(x_t|y)} \xrightarrow{\text{添加后条件}} = p(x_{t+1}|x_t) e^{\log p(y|x_{t+1}) - \log p(y|x_t)}$

$\frac{p(x_{t+1}|y)}{p(x_t|y)}$ 近似后可得

$$p(x_{t+1}|y) = N(x_{t+1}; \mu(x_t) + b_t^2 \nabla_x \log p(y|x_t), b_t^2 I)$$

$$\text{采样: } x_{t+1} = \mu(x_t) + \sigma_t^2 \nabla_{x_t} \log p(y|x_t) \Big|_{x_t=\mu(x_t)} + \epsilon_t$$

梯度缩放: 由分类器梯度加一个缩放参数 y , $\downarrow e^{y^2} \nabla_{x_t} \log p(y|x_t)$, 当 $y > 1$ 时, 会提高结果与 y 的相关性

② 由类 → 空间: $p(x_{t+1}|x_t, y) = \frac{p(x_{t+1}|x_t)}{Z_t} e^{y \sin(\theta_{x_t}, y)}$, Z_t 是归一化系数.

$$\xrightarrow{\text{近似}} p(x_{t+1}|x_t, y) = N(x_{t+1}; \mu(x_t) + \sigma_t^2 y \nabla_{x_t} \sin(\theta_{x_t}, y), \sigma_t^2 I)$$

$\stackrel{=0?}{\downarrow \text{SDE}}$ 余弦等相似度

③ 连续视角: SDE: 反向

$$dx = (f_t(x) - \frac{1}{2}(\sigma_t^2 + \sigma_t^2) \nabla_x \log p(x)) dt + \sigma_t dw$$

$$\therefore \nabla_x \log p_t(x_t|y) = \frac{\varepsilon_\theta(x_t, t) - \bar{\beta}_t \nabla_x \log p_t(y|x)}{\bar{\beta}_t} = \underbrace{\nabla_x \log p(x)}_{\varepsilon_\theta(x_t, t)} + \underbrace{\nabla_x \log p(y|x)}_{\bar{\beta}_t}$$

只需将 $\varepsilon_\theta(x_t, t) \rightarrow \varepsilon_\theta(x_t, t) - \bar{\beta}_t \nabla_x \log p_t(y|x)$

2. Classifier-free: 定义 $p(x_{t+1}|x_t, y) = N(x_{t+1}; \mu(x_t, y), \sigma_t^2 I)$

$$\rightarrow \mu(x_t, y) = \frac{1}{dt} |x_t - \frac{\sigma_t^2}{\bar{\beta}_t} \varepsilon_\theta(x_t, y, t)|$$

$$\rightarrow L = E_{x_0, y \sim p(x_0, y), \varepsilon \sim N(0, 1)} [\|\varepsilon - \varepsilon_\theta(x_t, y, t)\|^2]$$

为了平衡相关性和多样性, 令 $w = y^{-1}$

$$\hat{\varepsilon}_\theta(x_t, y, t) = (Hw) \varepsilon_\theta(x_t, y, t) - w \varepsilon_\theta(x_t, t)$$

$\hat{\varepsilon}_\theta(x_t, y, t) = (Hw) \varepsilon_\theta(x_t, y, t) - w \varepsilon_\theta(x_t, \phi, t)$ 引入一个特定输入 ϕ , 对应图像为全体图像 (条件)