



a) Diffusion Process

For input data: $x_0 \sim q(x_0)$

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t I)$$

$\{\beta_t\}_{t=1}^T$ is variance & $\beta_i \in [0,1]$ & $\beta_1 < \beta_2 < \dots < \beta_T$ (variance schedule)

if the T is huge enough, then $x_T \sim N(0,1)$. For every step:

$$q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t-1}) \quad (\text{Markov chain})$$

And the process is fixed, we need to define a variance schedule.
(eg. DDPM - linear)

Def: $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$. Through reparameterization trick

$$x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon_{t-1}, \text{ where } \epsilon_{t-1}, \epsilon_{t-2}, \dots \sim N(0, I)$$

$$= \sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} \epsilon_{t-1}$$

$$= \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_{t-1}} \epsilon_{t-2}) + \sqrt{1-\alpha_t} \epsilon_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{(\alpha_t - \alpha_t \alpha_{t-1}) + (1-\alpha_t)^2} \bar{\epsilon}_{t-2}, \text{ where } \bar{\epsilon}_{t-2} \text{ merges two Gaussians (*)}$$

$$= \dots$$

$$= \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon$$

$$\begin{aligned} &\Downarrow \\ &\epsilon_1 \sim N(0, \sigma_1^2 I), \epsilon_2 \sim N(0, \sigma_2^2 I) \\ &\oplus \\ &\epsilon \sim N(0, (\sigma_1^2 + \sigma_2^2) I) \end{aligned}$$

through variational reparameterization:

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

$\sqrt{\bar{\alpha}_t}$: signal - ~~noise~~ rate

$\sqrt{1-\bar{\alpha}_t}$: noise - rate

② b) reverse denoising process (generative)

define the reverse denoising process as a Markov chain. by Networks.
To calculate $q(x_{t+1}|x_t)$, we could use a parameterized Gaussian distribution

$$P_\theta(x_{0:T}) = p(x_T) \prod_{i=1}^T p_\theta(x_{t-1}|x_t)$$

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

where $p(x_T) = \mathcal{N}(x_T; 0, I)$

the $p_\theta(x_{t-1}|x_t)$ is parameterized Gaussian distribution, their mean & variance are got by training networks $\mu_\theta(x_t, t), \Sigma_\theta(x_t, t)$. In fact, the diffusion model aims to get these trained networks cause they are part of the final diffusion model.

Though the $q(x_t|x_{t+1})$ can't be calculate, the prior distribution $q(x_{t+1}|x_t, x_0)$ could be done. There's

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

1. Bayesian condition equation:

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

2. Cause diffusion process Markov chain, we know that

$$q(x_t|x_{t+1}, x_0) = q(x_t|x_{t+1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t} x_{t+1}, \beta_t I)$$

and,

$$q(x_{t+1}|x_0) = \mathcal{N}(x_{t+1}; \sqrt{\alpha_{t+1}} x_0, (1-\alpha_{t+1}) I),$$

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

so,

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

$$\left(-\frac{1}{2} \left(\frac{(x_t - \sqrt{\alpha_t} x_{t+1})^2}{\beta_t} + \frac{(x_{t+1} - \sqrt{\alpha_{t+1}} x_0)^2}{1-\alpha_{t+1}} - \frac{(x_t - \sqrt{\alpha_t} x_0)^2}{1-\alpha_t} \right) \right)$$

$$= \exp \left(-\frac{1}{2} \left(\frac{x_t^2 - 2\sqrt{\alpha_t} x_t x_{t+1} + \alpha_t x_{t+1}^2}{\beta_t} + \frac{x_{t+1}^2 - 2\sqrt{\alpha_{t+1}} x_{t+1} x_0 + \alpha_{t+1} x_0^2}{1-\alpha_{t+1}} - \frac{x_t^2 - 2\sqrt{\alpha_t} x_t x_0 + \alpha_t x_0^2}{1-\alpha_t} \right) \right)$$

$$= \exp \left(-\frac{1}{2} \left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\alpha_{t+1}} \right) x_{t+1}^2 - \left(\frac{2\sqrt{\alpha_t}}{\beta_t} x_t + \frac{2\sqrt{\alpha_{t+1}}}{1-\alpha_{t+1}} x_0 \right) x_{t+1} + C(x_t, x_0) \right) \right)$$

There $C(x_t, x_0)$ has no relation with ~~x_{t-1}~~ x_{t-1} .

According to Gaussian distribution definitions we could get prior distribution ^{mean} & ^{variance}.

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

$$\begin{aligned} \tilde{\mu}_t(x_t, x_0) &= \left(\frac{\sqrt{\alpha_t}}{\beta_t} x_t + \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \\ &= \left(\frac{\sqrt{\alpha_t}}{\beta_t} x_t + \frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) \cdot \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{1 - \bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} x_0 \end{aligned}$$

we could see that $\hat{\beta}_t$ is constant, and the $\tilde{\mu}_t$ is function rely on x_t & x_0 .

c). optimization

if we consider the x_t as a latent variable, and the diffusion model could be a hierarchical VAEs. so we could get ELBO (variational lower bound) through variational inference. there

$$\log p_\theta(x_0) = \log \int p_\theta(x_{0:T}) dx_{1:T}$$

$$= \log \int \frac{p_\theta(x_{0:T}) q(x_{1:T}|x_0)}{q(x_{1:T}|x_0)} dx_{1:T} \geq E_{q(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right]$$

↓

$$-\log p_\theta(x_0) \leq -\log p_\theta(x_0) + D_{KL}(q(x_{1:T}|x_0) || p_\theta(x_{1:T}|x_0))$$

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

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

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

For networks training, we put negative in front of VLB:

$$\mathcal{L} = -\mathcal{L}_{VLB} = E_{q(x_{1:T}|x_0)} \left[-\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right]$$

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

④

and then:

$$\begin{aligned}
L &= E_{q(X_{1:T}|X_0)} \left[\log \frac{q(X_{1:T}|X_0)}{p_\theta(X_{0:T})} \right] = E_q \left[\log \frac{\prod_{t=1}^T q(X_t|X_{t-1})}{p_\theta(X_T) \prod_{t=1}^T p_\theta(X_{t-1}|X_t)} \right] \\
&= E_q \left[-\log p_\theta(X_T) + \sum_{t=1}^T \log \frac{q(X_t|X_{t-1})}{p_\theta(X_{t-1}|X_t)} \right] \\
&= E_q \left[-\log p_\theta(X_T) + \sum_{t=2}^T \log \frac{q(X_t|X_{t-1})}{p_\theta(X_{t-1}|X_t)} + \log \frac{q(X_1|X_0)}{p_\theta(X_0|X_1)} \right] \\
&= E_q \left[-\log p_\theta(X_T) + \sum_{t=2}^T \log \frac{q(X_t|X_{t-1}, X_0)}{p_\theta(X_{t-1}|X_t)} + \log \frac{q(X_1|X_0)}{p_\theta(X_0|X_1)} \right] \quad (\text{add prior } X_0) \\
&= E_q \left[-\log p_\theta(X_T) + \sum_{t=2}^T \log \left[\frac{q(X_{t-1}|X_t, X_0)}{p_\theta(X_{t-1}|X_t)} \cdot \frac{q(X_t|X_0)}{q(X_{t-1}|X_0)} \right] + \log \frac{q(X_1|X_0)}{p_\theta(X_0|X_1)} \right] \\
&\quad \text{(Bayesian method)} \\
&= E_q \left[-\log p_\theta(X_T) + \sum_{t=2}^T \log \frac{q(X_{t-1}|X_t, X_0)}{p_\theta(X_{t-1}|X_t)} + \log \frac{q(X_T|X_0)}{q(X_1|X_0)} + \log \frac{q(X_1|X_0)}{p_\theta(X_0|X_1)} \right] \\
&= E_q \left[\log \frac{q(X_1|X_0)}{p_\theta(X_T)} + \sum_{t=2}^T \log \frac{q(X_{t-1}|X_t, X_0)}{p_\theta(X_{t-1}|X_t)} - \log p_\theta(X_0|X_1) \right] \\
&= E_q \left[\log \frac{q(X_1|X_0)}{p_\theta(X_T)} \right] + \sum_{t=2}^T E_{q(X_t|X_0)} \left[\log \frac{q(X_{t-1}|X_t, X_0)}{p_\theta(X_{t-1}|X_t)} - E_{q(X_1|X_0)} [\log p_\theta(X_0|X_1)] \right] \\
&= E_{q(X_T|X_0)} \left[\log \frac{q(X_1|X_0)}{p_\theta(X_T)} \right] + \sum_{t=2}^T E_{q(X_t|X_0)} \left[\log \frac{q(X_{t-1}|X_t, X_0)}{p_\theta(X_{t-1}|X_t)} \right] - E_{q(X_1|X_0)} [\log p_\theta(X_0|X_1)] \\
&= \underbrace{D_{KL}(q(X_T|X_0) || p_\theta(X_T))}_{L_T} + \underbrace{\sum_{t=2}^T E_{q(X_t|X_0)} [D_{KL}(q(X_{t-1}|X_t, X_0) || p_\theta(X_{t-1}|X_t))]}_{L_{t-1}} - \underbrace{E_{q(X_1|X_0)} [\log p_\theta(X_0|X_1)]}_{L_0}
\end{aligned}$$

L_0 : ~~the~~ rebuild original data:

$$p_\theta(X_0|X_1) = \prod_{i=1}^D \int_{\delta_-(x_0^i)}^{\delta_+(x_0^i)} \mathcal{N}(x_0; \mu_\theta^i(x_1, 1), \Sigma_\theta^i(x_1, 1)) dx$$

$$\delta_+(x) = \begin{cases} \infty & x=1 \\ x + \frac{1}{255} & x < 1 \end{cases}, \quad \delta_-(x) = \begin{cases} -\infty & x=-1 \\ x - \frac{1}{255} & x > -1 \end{cases}$$

⑤ ^{the KL converge between}
 L_T : to calculate noise distribution and prior distribution. parameterization
 $p(X_T) = \mathcal{N}(0, I)$, $q(X_T | X_0) = \mathcal{N}(0, I)$, the $L_T \rightarrow 0$ with non trained-

L_{t-1} : to calculate the KL between $p_\theta(X_{t+1} | X_t)$ and $q(X_{t+1} | X_t, X_0)$. The $q(X_{t+1} | X_t, X_0)$ is a Gaussian distribution, and that's why we def $p_\theta(X_{t+1} | X_t)$ as a Gaussian distribution by networks parameterization.

For both L_0 & L_{t-1} , we want to get trained-well networks $\mu_\theta(X_t, t)$ & $\Sigma_\theta(X_t, t)$ (to $L_0, t=1$). Here, DDPM simplified the object:

$\Sigma_\theta(X_t, t) = \sigma_t^2 I$. ($\sigma_t^2 = \beta_t$ or $\tilde{\beta}_t$, or trainable variance;
 there we suppose $\sigma_t^2 = \tilde{\beta}_t$. so,

$$q(X_{t+1} | X_t, X_0) = \mathcal{N}(X_{t+1}; \tilde{\mu}(X_t, X_0), \sigma_t^2 I)$$

$$p_\theta(X_{t+1} | X_t, X_0) = \mathcal{N}(X_{t+1}; \mu_\theta(X_t, t), \sigma_t^2 I)$$

To calculate the KL ~~betw~~ with two Gaussian distribution:

$$KL(p_1 || p_2) = \frac{1}{2} \left[\text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_1^{-1} (\mu_2 - \mu_1) - n + \log \frac{\det(\Sigma_2)}{\det(\Sigma_1)} \right]$$

so, we have

$$\begin{aligned} D_{KL}[q(X_{t+1} | X_t, X_0) || p_\theta(X_{t+1} | X_t)] &= D_{KL}[\mathcal{N}(X_{t+1}; \tilde{\mu}(X_t, X_0), \sigma_t^2 I) || \mathcal{N}(X_{t+1}; \mu_\theta(X_t, t), \sigma_t^2 I)] \\ &= \frac{1}{2} \left[n + \frac{1}{\sigma_t^2} \|\mu_\theta(X_t, t) - \tilde{\mu}(X_t, X_0)\|^2 - n + \log 1 \right] \\ &= \frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(X_t, t) - \mu_\theta(X_t, t)\|^2. \end{aligned}$$

and, L_{t-1} :

$$L_{t-1} = E_{q(X_t | X_0)} \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}(X_t, t) - \mu_\theta(X_t, t)\|^2 \right]$$

According to the equation, we want to get the difference between $\tilde{\mu}(X_t, t)$ and $\mu_\theta(X_t, t)$ become smaller.

⑥ But DDPM found that if we predict ε_0 rather than $\mu_0(x_t, t)$ it could be better.
According to diffusion process:

$$x_t(x_0, \varepsilon) = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \varepsilon, \text{ where } \varepsilon \sim \mathcal{N}(0, 1)$$

put it into L_{t-1} :

$$\begin{aligned} L_{t-1} &= E_{x_0} \left(E_{\varepsilon(x_t|x_0)} \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_0(x_t, t)\|^2 \right] \right) \\ &= E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(x_t(x_0, \varepsilon), \frac{1}{\sqrt{\alpha_t}} (x_t(x_0, \varepsilon) - \sqrt{1-\alpha_t} \varepsilon)) - \mu_0(x_t(x_0, \varepsilon), t)\|^2 \right] \\ &= E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\frac{1}{2\sigma_t^2} \left\| \left(\frac{\sqrt{\alpha_t}(1-\alpha_{t-1})}{1-\alpha_t} x_t(x_0, \varepsilon) + \frac{\sqrt{\alpha_{t-1}}\beta_t}{1-\alpha_t} \cdot \frac{1}{\sqrt{\alpha_t}} (x_t(x_0, \varepsilon) - \sqrt{1-\alpha_t} \varepsilon) \right) - \mu_0(x_t(x_0, \varepsilon), t) \right\|^2 \right] \\ &= E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} (x_t(x_0, \varepsilon) - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon) - \mu_0(x_t(x_0, \varepsilon), t) \right\|^2 \right]. \end{aligned}$$

then, we reparameterized $\mu_0(x_t(x_0, \varepsilon), t)$:

$$\mu_0(x_t(x_0, \varepsilon), t) = \frac{1}{\sqrt{\alpha_t}} \left[x_t(x_0, \varepsilon) - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon_0(x_t(x_0, \varepsilon), t) \right]$$

there ε_0 is a trained function by network. That means we ~~use ε_0~~ ^{aims to} predict ε rather than $\mu_0(x_t, t)$. Then,

$$\begin{aligned} L_{t-1} &= E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} (x_t(x_0, \varepsilon) - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon) - \frac{1}{\sqrt{\alpha_t}} \left(x_t(x_0, \varepsilon) - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon_0 \right) \right\|^2 \right] \\ &= E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1-\alpha_t)} \|\varepsilon - \varepsilon_0(x_t(x_0, \varepsilon), t)\|^2 \right] \\ &= E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1-\alpha_t)} \|\varepsilon - \varepsilon_0(\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon, t)\|^2 \right] \end{aligned}$$

DDPM simplified training object function:

$$L_{t-1}^{\text{simple}} = E_{x_0, \varepsilon \sim \mathcal{N}(0, 1)} \left[\|\varepsilon - \varepsilon_0(\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon, t)\|^2 \right]$$