

DDPM loss design

We now know q and p are both normal distribution of the probability projection model (diffusion model)

We proved the each step of denoising can be numerically calculated

We want the model's parameters (θ) can be better, by having a better estimation of each time steps' denoising process. (When the model defined by θ can have a closer denoising to the theoretical 'reverse diffuse')

We use the KL divergence to measure the difference of the processes, and we use the maximum likelihood estimation to seek the θ parameters. So we can build a loss that supervising each step t .

\Rightarrow In Diffusion we have

$$\begin{cases} q(x_t | x_{t-1}) & \overset{t-1}{0} \rightarrow \overset{t}{0} & \text{diffuse} \\ q(x_{t-1} | x_t) & \overset{t-1}{0} \leftarrow \overset{t}{0} & \text{reverse diffuse} \\ p_\theta(x_{t-1} | x_t) & \overset{t-1}{0} \leftarrow \overset{t}{0} & \text{denoise} \end{cases}$$

and $q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$ (all the diffuse)
 $p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$ (all the denoise)

Our aim is to use p_θ to estimate the reverse q .
 if we can have all the x_0, x_1, \dots, x_T

the process supervise Target here is to find a θ

that the $p_\theta(x_{0:T} | x_0)$ is very close to $q(x_{0:T} | x_0)$

So we need two distribution (abstractive) things are very close

Measure KL Divergence (information theory)

Def:

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} P(x) \log \frac{P(x)}{Q(x)} dx \quad (\text{Def})$$

As we know

$$1. D_{KL}(P||Q) \neq D_{KL}(Q||P) \quad (0)$$

$$2. D_{KL}(P||Q) \geq 0, "=" \text{ at } P=Q \quad (1)$$

Maximum likelihood estimation (MLE)

how we know we can calculate $\mu_\theta(x_t, t)$ and $\Sigma_\theta(x_t, t)$ of p_θ
 but how to get the correct θ so the value works?

MLE: L is a function about θ ; when we have the given output of X , the value of L is possibility of obtain X under θ :

$$L(\theta|X) = P(X=x|\theta) \quad \text{we want the chance to be the highest (maximum-likelihood)}$$

$$\Rightarrow \min -L \Rightarrow \min -L(\theta|X) = \min -P(X=x|\theta)$$

How to build this L ? $P_\theta(x_0) := \int P_\theta(x_{0:T}) dx_{1:T}$ by right this is the mathematic target def.

$$L = E_q[-\log P_\theta(x_0)]$$

We have $D_{KL}(q(x_{1:T}|x_0) || P_\theta(x_{1:T}|x_0)) \geq 0$ (Def) (1)

$$\Rightarrow -\log P_\theta(x_0) \leq -\log P_\theta(x_0) + D_{KL}(q(x_{1:T}|x_0) || P_\theta(x_{1:T}|x_0)) \quad \text{Def?}$$

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

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

$$\Rightarrow -\log P_\theta(x_0) \leq E_q \left[\log \frac{q(x_{1:T}|x_0)}{P_\theta(x_{0:T})} \right] \quad \begin{array}{l} \text{Variational lower bound (VLB)} \\ \text{Evidence lower bound (ELBO)} \end{array}$$

$$\text{let } L_{VLB} = E_{q(x_{0:T})} \left[\log \frac{q(x_{1:T}|x_0)}{P_\theta(x_{0:T})} \right] \geq -E_{q(x_0)} \log P_\theta(x_0)$$

Then

$$\begin{aligned} L_{VLB} &= E_{q(x_{0:T})} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right] \\ &= E_q \left[\log \frac{\prod_{t=1}^T q(x_t|x_{t-1})}{p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)} \right] \\ &= E_q \left[-\log p_\theta(x_T) + \sum_{t=1}^T \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)} \right] \\ &= E_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right] \\ &= E_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} \cdot \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} \right) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right] \\ &= E_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \sum_{t=2}^T \log \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
&= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \\
&= \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right]
\end{aligned}$$

$$\Rightarrow \mathcal{L}_{\text{VLB}} := L_T + L_{T-1} + \dots L_0$$

$$L_T := D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))$$

$$L_{t-1} := D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)), \quad 1 \leq t \leq T$$

$$L_0 := -\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$$

for L_T : its based on $\mathbf{x}_T \sim N(0,1)$.

So we can ignore it. (first step can go anywhere)
 why? q has no parameter. and p_θ cannot be supervised based on $\mathbf{x}_T \sim N(0,1)$

for L_0 : not very helpful.

the prove will be explored later.

for L_t : $D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)), \quad 1 \leq t \leq T-1$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ can be numerically solve

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_t \right)$$

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

dim dim dim
 mesh mesh mesh

$p_0(x_{t+1}|x_t)$ is also a gaussian distribution (move + rescale)

given as $p_0(x_{t+1}|x_t) = N(x_{t+1}; \mu_0(x_t, t), \Sigma_0(x_t, t))$

the D_{KL} of two gaussian p, q can be given

$$D_{KL}(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$



to optimize the D_{KL} of p_0 and q .

the $\tilde{\beta}_t$ of q and $\Sigma_0(x_t, t)$ of p are all constant which are invariance to θ optimization. (remove)

\Rightarrow So \mathcal{L}_t only consider about $\tilde{\mu}_t$ and $\mu_0(x_t, t)$

$$\mathcal{L}_t = E_q[\|\tilde{\mu}_t(x_t, x_0) - \mu_0(x_t, t)\|^2]$$

$$= E_{x_0, \epsilon}[\|\frac{1}{\sqrt{\alpha_t}}(\underbrace{x_t(x_0, \epsilon) - \frac{\beta_t}{\sqrt{1-\alpha_t}}\epsilon}_{\epsilon \sim N(0,1)} - \mu_0(x_t(x_0, \epsilon), t))\|^2]$$

both of the mean value is about x_t under x_0 , and ϵ



assume p_0 and q (reverse) has the same mean value

$$\Rightarrow \mu_0(x_t(x_0, \epsilon), t) = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \underbrace{\epsilon_0(x_t, t)}_{\text{unknown}})$$

put it into the above \mathcal{L}_t , to get the lowest D_{KL} is now get two means that are close to each other

$$\mathcal{L}_t = E_{x_0, \epsilon}[\|\frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \underbrace{\epsilon}_{\text{noise estimation}}) - \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \underbrace{\epsilon_0(x_t, t)}_{\text{noise estimation}})\|^2]$$

mean value estimation now become noise estimation, $\epsilon \sim N(0,1)$

$$\mathcal{L}_t \propto E_{X_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2], \quad \epsilon \sim N(0, 1)$$

ϵ_θ is a noise given by x_t and timestep t . remove all constant

$$= E_{X_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2], \quad \epsilon \sim N(0, 1)$$

$$\mathcal{L}_t \Rightarrow l_2\text{-loss} \rightarrow \begin{matrix} \epsilon_\theta(\sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \\ \epsilon, \epsilon \sim N(0, 1) \end{matrix}$$

loss is define \uparrow to the l_2 -loss at the timestep $t \in [0, T]$

$$\text{loss}_{\text{simple}}(\theta) := E_{t, X_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2]$$

$\text{dimer} \rightarrow \text{image} + \epsilon \cdot \text{something}$

$$\Rightarrow \mathcal{L}_{\text{simple}} = E_{t, X_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2]$$

estimating a noise ϵ ,
 who was used to generate a dimer & noised "Image" at t .
 based on it

$$\text{loss} = l_2(\epsilon, I(\epsilon, t, \beta))$$

$\uparrow \quad \uparrow \quad \uparrow$
 time step all settings