# 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 (theta) can be better, by having a better estimation of each time steps' denoising process. (When the model defined by theta 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 theta parameters. So we can build a loss that supervising each step t.

$\Rightarrow$ in Diffusion $\qquad \Big\{ \begin{array}{l} q\,(x_t \,|\, x_{t-1}) \quad \overset{t-1}{0} \to \overset{t}{0} \quad \text{diffuse} \\ q\,(x_{t-1} \,|\, x_t) \quad \overset{t-1}{0} \leftarrow \overset{t}{0} \quad \text{reverse diffuse} \\ P_\theta\,(x_{t-1} \,|\, x_t) \quad \overset{t-1}{0} \leftarrow \overset{t}{0} \quad \text{denoise}\ \theta \end{array}$

**Previous**  we have ✕

and, $\quad q\,(x_{1:T} \,|\, x_0) = \prod\limits_{t=1}^{T} q\,(x_t \,|\, x_{t-1}) \qquad$ (all the diffuse)

$2\ P_\theta\,(x_{0:T}) = P(x_T)\, \prod\limits^{T} P_\theta\,(x_{t-1} \,|\, x_t) \quad$ (all the denoise)

$\mathcal{O}$ur aim is to use $P_\theta$ to estimate the reverse $q$.

if we can have all the $x_0, X_1 \cdots X_T$ $\qquad \begin{array}{l} 0 \to 0 \to 0 \to \overset{X_1}{0} \\ X_0 \leftarrow \leftarrow \hookleftarrow \end{array}$

the process **supervise** Target here is to find a $\theta$

who makes $P_\theta\,(x_{0:T} \,|\, x_0)$ 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:
$$D_{KL}\,(P \| Q) = \int_{-\infty}^{\infty} P(x)\, \log \frac{P(x)}{q(x)}\, dx \qquad \text{(Def)}$$

As we know

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

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

**Maximum likelyhood estimation** (MLE) $\quad \swarrow$ Variance

how we know we can calculate $M_\theta\,(x_t, t)$ and $\Sigma_\theta\,(x_t, t)$ of $P_\theta$

but how to get the correct $\theta$ so the value works?

denoise

MLE:  $\mathcal{L}$ is a function about $\Theta$ ; when we have the given

output of $x$ . the value of $\mathcal{L}$ is possibility of obtain $x$ under $\Theta$:

$$\mathcal{L}(\Theta|x) = P(X=x|\Theta)$$  we want the chance to be the highest
(maximumazation)

also as

$\Rightarrow$ min $-\mathcal{L}$ $\Rightarrow$ min $-\mathcal{L}(\Theta|x) = $ min $-P(X=x|\Theta)$ $P_\Theta(X)$

How to build this $\mathcal{L}$ ?

$$P_\Theta(x_0) := \int P_\Theta(x_{0:T}) dx_{1:T}$$  by right this is the mathmatic Target def.

Search for the minimal negative log Expectation

$$\mathcal{L} = E_q(-\log P_\Theta(x_0))$$ 最小化负对数 似然  $0 \to 0 \to 0$
                                                    $x_0$    $T$

We have

$$D_{KL}\left(q(x_{1:T}|x_0) \| P_\Theta(x_{1:T}|x_0)\right) \geq 0 \qquad (Def) \ [1]$$

add something $\geq 0$

$\Rightarrow$  $-\log P_\Theta(x_0) \leq -\log P_\Theta(x_0) + D_{KL}\left(q(x_{1:T}|x_0)\| P_\Theta(x_{1:T}|x_0)\right)$  (Def?)

$$E\left[-\log P_\Theta(x_0)\right] \leq E\left[-\log P_\Theta(x_0)\right] + E_{x_{1:T} \sim q(x_{1:T}|x_0)}\left[\log \frac{q(x_{1:T}|x_0)}{P_\Theta(x_{0:T})/P_\Theta(x_0)}\right]$$

remove

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

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

Variational lower bound (VLB)
Evidence lower bound (ELBO)
VAE

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

Then
$$L_{\text{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})}\left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})}\right]$$

$$= \mathbb{E}_q\left[\log \frac{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T)\prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}\right]$$

$$= \mathbb{E}_q\left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}\right]$$

$$= \mathbb{E}_q\left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}\right]$$

$$= \mathbb{E}_q\left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \left(\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \cdot \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}\right) + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}\right]$$

$$= \mathbb{E}_q\left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}\right]$$

$$= \mathbb{E}_q\left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^t \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}\right]$$

$$= \mathbb{E}_q\left[ \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)\right]$$

$$= \mathbb{E}_q[\underbrace{D_{\mathrm{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \,\|\, p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\mathrm{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \,\|\, p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} \underbrace{- \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0}]$$

-)

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

$$L_T := D_{KL}\left( q(X_T|X_0) \| P_\theta(X_T) \right)$$

$$L_{t-1} := D_{KL}\left( q(X_{t-1}|(X_t, X_0)) \| P_\theta(X_{t-1}|X_t) \right), \quad 1 \leq t \leq T$$

$$L_0 := -\log P_\theta(X_0|X_1)$$

for $L_T$ : its based on $X_T \sim N(0,1)$.

So we can ignore it. (first step can go anywhere)

why? $q$ has no parameter. and $P_\theta$ cannot be supervised based on $X_T \sim N(0,1)$

for $L_0$ :

not very helpful.

→ the prove will be explored lavler.

for $L_t$: $D_{KL}\left( q(X_{t-1}|(X_t, X_0)) \| P_\theta(X_{t-1}|X_t) \right), \quad 1 \leq t \leq T-1$

$q(X_{t-1}|(X_t, X_0))$ can be numerically solve

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}}\left(X_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_t\right)$$

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

dim dim dim
mesh mesh mesh

$P_\theta(x_{t-1}|x_t)$ is also a gaussian distribution ( move + rescale)

given as $P_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t,t), \Sigma_\theta(x_t,t))$

---

the $D_{KL}$ of two gaussian $p,q$ $\overset{(\mu_1,\sigma_1)\ (\mu_2,\sigma_2)}{\text{can be given}}$  Formular. of $D_{KL}(p,q)$

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

---

to optimize the $D_{KL}$ of $P_\theta$ and $q$.

the $\tilde{\beta}_t$ of $q$ and $\Sigma_\theta(x_t,t)$ of $p$ are all constant which are invariance to $\theta$ optimization. (remove) in $\mathcal{L}$

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

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

$$= E_{x_0,\varepsilon}\left[\left\| \frac{1}{\sqrt{\alpha_t}}\left(x_t(x_0,\varepsilon) - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon\right) - \mu_\theta(x_t(x_0,\varepsilon),t) \right\|^2\right]$$

both of the mean value is about $x_t$ under $x_0$, and $\varepsilon$    $\varepsilon \sim N(0,1)$

assume $P_\theta$ and $q$ (reverse) has the same mean value

=)

$$\mu_\theta(x_t(x_0,\varepsilon),t) = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_\theta(x_t,t)\right)$$
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

$$\mathcal{L}_t = E_{x_0,\varepsilon}\left[\left\| \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon\right) - \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_\theta(x_t,t)\right) \right\|^2\right]$$

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

$$L_t \propto E_{X_0, \varepsilon}\left[\|\varepsilon - \varepsilon_\theta(X_t, t)\|^2\right], \quad \varepsilon \sim N(0,1) \quad \text{remove all constant}$$

$\varepsilon_\theta$ is a noise given by $X_t$ and time step $t$.

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

$$L_t \Rightarrow l_2 - loss \rightarrow \varepsilon_\theta(\sqrt{\bar{\alpha}_t}X_0 + \sqrt{1-\bar{\alpha}_t}\varepsilon, t)$$

$$\searrow \varepsilon, \quad \varepsilon \sim N(0,1)$$

---

loss $\uparrow$ to is define

the $l_2$-loss at the time step $t \in [0, T]$

$$loss_{simple}(\theta) := E_{t, X_0, \varepsilon}\left[\|\varepsilon - \varepsilon_\theta(\underline{\sqrt{\bar{\alpha}_t}X_0 + \sqrt{1-\bar{\alpha}_t}\varepsilon}, t)\|^2\right]$$

diner image + $\varepsilon \cdot$ something

$$\Rightarrow L_{simple} = E_{t, X_0, \varepsilon}\left[\|\varepsilon - \varepsilon_\theta(X_t, t)\|^2\right]$$

estimating a noise $\varepsilon$,
who was used to generate a diner & noised "Image" at $t$.
based on it

$$loss = l_2\left(\varepsilon, I(\varepsilon, t, \beta)\right)$$

time step   all settings

At training and inference time, we know the **β**'s, **α**'s, and $\mathbf{x_t}$. So our **model only needs to predict the noise at each timestep.** The simplified (after ignoring some weighting terms) loss function used in the *Denoising Diffusion Probabilistic Models* is as follows:

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t,\mathbf{x}_0,\epsilon}\left[\left\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t)\right\|^2\right]$$

*Comparing just the noise.*

Which is basically:

$$L_{\text{simple}} = \mathbb{E}_{t,x_0,\epsilon}\left[\|\epsilon - \epsilon_\theta(x_t, t)\|^2\right]$$

***This is the final loss function we use to train DDPMs, which is just a "Mean Squared Error" between the noise added in the forward process and the noise predicted by the model. This is the most impactful contribution of the paper Denoising Diffusion Probabilistic Models.***

It's awesome because, beginning from those scary-looking ELBO terms, we ended up with the simplest loss function in the entire machine learning domain.