# Variational Inference for Diffusion Models

Tomonari MASADA

masada@rikkyo.ac.jp

# Contents

# Marginal likelihood

We consider a Bayesian generative model whose joint distribution is

$$p_\theta(\mathbf{x}_{0:T}) = p_\theta(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \tag{1}$$

where $\mathbf{x}_0$ is an observation and $\{\mathbf{x}_t : t = 1, \ldots, T\}$ are latent random variables. The distribution of $\mathbf{x}_t$ is only conditioned on $\mathbf{x}_{t+1}$. The log of the marginal likelihood ($=$ the evidence) of $\mathbf{x}_0$ is

$$\log p_\theta(\mathbf{x}_0) = \log \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \tag{2}$$

# Contents

# ELBO

Jensen's inequality gives a lower bound of $\log p_\theta(\mathbf{x}_0)$ as follows:

$$
\begin{aligned}
\log p_\theta(\mathbf{x}_0) &= \log \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \frac{p_\theta(\mathbf{x}_{0:T})}{q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&\geq \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{0:T})}{q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{0:T})}{q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \equiv L_{\text{VLB}} \quad (3)
\end{aligned}
$$

where $q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0)$ is a variational posterior, which is conditioned on the observation $\mathbf{x}_0$ as in VAE. We train our model over many observations $\mathcal{X} \equiv \{\mathbf{x}_0^{(1)}, \ldots, \mathbf{x}_0^{(N)}\}$ in an amortized manner.

# Markov assumption

The definition of the conditional distribution gives the following equations:

$$q_\psi(\mathbf{x}_2|\mathbf{x}_1, \mathbf{x}_0)q_\psi(\mathbf{x}_1|\mathbf{x}_0) = \frac{q_\psi(\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_0)}{q_\psi(\mathbf{x}_1, \mathbf{x}_0)}\frac{q_\psi(\mathbf{x}_1, \mathbf{x}_0)}{q_\psi(\mathbf{x}_0)} = \frac{q_\psi(\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_0)}{q_\psi(\mathbf{x}_0)}$$

$$= q_\psi(\mathbf{x}_2, \mathbf{x}_1|\mathbf{x}_0) \tag{4}$$

$$q_\psi(\mathbf{x}_3|\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_0)q_\psi(\mathbf{x}_2|\mathbf{x}_1, \mathbf{x}_0)q_\psi(\mathbf{x}_1|\mathbf{x}_0) = \frac{q_\psi(\mathbf{x}_3, \mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_0)}{q_\psi(\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_0)}\frac{q_\psi(\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_0)}{q_\psi(\mathbf{x}_1, \mathbf{x}_0)}\frac{q_\psi(\mathbf{x}_1, \mathbf{x}_0)}{q_\psi(\mathbf{x}_0)}$$

$$= q_\psi(\mathbf{x}_3, \mathbf{x}_2, \mathbf{x}_1|\mathbf{x}_0) \tag{5}$$

$$\cdots$$

Therefore, by assuming that the equation $q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \ldots, \mathbf{x}_1, \mathbf{x}_0) = q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$ holds for $t = 2, \ldots, T$, we obtain the following factorization of $q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0)$:

$$q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) = q_\psi(\mathbf{x}_1|\mathbf{x}_0)\prod_{t=2}^{T} q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \tag{6}$$

Then the variational lower bound $L_{\text{VLB}}$ can be rewritten as follows:

$$
\begin{aligned}
L_{\text{VLB}} &= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{0:T})}{q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^{T} q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log p_\theta(\mathbf{x}_T) d\mathbf{x}_{1:T} + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \sum_{t=2}^{T} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&\quad + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q_\psi(\mathbf{x}_1|\mathbf{x}_0)} d\mathbf{x}_{1:T}
\end{aligned}
\tag{7}
$$

Using Bayes' rule, we have

$$q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q_\psi(\mathbf{x}_t|\mathbf{x}_0)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_0)} \tag{8}$$

We will discuss later how we make the above $q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$ tractable.

We replace $q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$ appearing in $L_{\text{VLB}}$ with $\frac{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q_\psi(\mathbf{x}_t|\mathbf{x}_0)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_0)}$ based on Eq. (8) and rewrite $L_{\text{VLB}}$ as follows:

$$L_{\text{VLB}} = \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log p_\theta(\mathbf{x}_T)d\mathbf{x}_{1:T} + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \sum_{t=2}^{T} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}d\mathbf{x}_{1:T}$$

$$+ \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \sum_{t=2}^{T} \log \frac{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q_\psi(\mathbf{x}_t|\mathbf{x}_0)}d\mathbf{x}_{1:T} + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q_\psi(\mathbf{x}_1|\mathbf{x}_0)}d\mathbf{x}_{1:T} \tag{9}$$

(continued on the next page)

$$
\begin{aligned}
L_{\text{VLB}} &= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log p_\theta(\mathbf{x}_T) d\mathbf{x}_{1:T} + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&\quad + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{q_\psi(\mathbf{x}_1|\mathbf{x}_0)\cancel{q_\psi(\mathbf{x}_2|\mathbf{x}_0)}\cdots\cancel{q_\psi(\mathbf{x}_{T-1}|\mathbf{x}_0)}}{\cancel{q_\psi(\mathbf{x}_2|\mathbf{x}_0)}\cancel{q_\psi(\mathbf{x}_2|\mathbf{x}_0)}\cdots q_\psi(\mathbf{x}_T|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&\quad + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q_\psi(\mathbf{x}_1|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log p_\theta(\mathbf{x}_T) d\mathbf{x}_{1:T} + \sum_{t=2}^T \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&\quad + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{\cancel{q_\psi(\mathbf{x}_1|\mathbf{x}_0)}}{q_\psi(\mathbf{x}_T|\mathbf{x}_0)} d\mathbf{x}_{1:T} + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{\cancel{q_\psi(\mathbf{x}_1|\mathbf{x}_0)}} d\mathbf{x}_{1:T} \\
&= \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_T)}{q_\psi(\mathbf{x}_T|\mathbf{x}_0)} d\mathbf{x}_{1:T} + \sum_{t=2}^T \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&\quad + \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) d\mathbf{x}_{1:T} \equiv L_T + \sum_{t=2}^T L_{t-1} + L_0 \quad\quad (10)
\end{aligned}
$$

We can rewrite $L_{t-1}$ as follows:

$$
\begin{aligned}
L_{t-1} &\equiv \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\
&= \int \Big( q_\psi(\mathbf{x}_1|\mathbf{x}_0) \prod_{t' \neq t} q_\psi(\mathbf{x}_{t'-1}|\mathbf{x}_{t'}, \mathbf{x}_0) \Big) \Big( q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \Big) d\mathbf{x}_{1:T} \\
&= -\int \Big( q_\psi(\mathbf{x}_1|\mathbf{x}_0) \prod_{t' \neq t} q_\psi(\mathbf{x}_{t'-1}|\mathbf{x}_{t'}, \mathbf{x}_0) \Big) D_{\mathsf{KL}}(q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) d\mathbf{x}_{1:T} \\
&\equiv -\mathbb{E}_{\neg t}\big[ D_{\mathsf{KL}}(q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) \big] \qquad (11)
\end{aligned}
$$

It can be said that, by minimizing this expectation of the KL-divergence
$D_{\mathsf{KL}}(q_\psi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$ for each $t$, we can maximize the ELBO in Eq. (7).

# Contents

# Parameterization of variational posterior

We parameterize our variational posterior with the parameters $\psi \equiv \{\alpha_t : t = 1, \ldots, T\}$ as

$$q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I}) \tag{12}$$

Therefore, $q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1})$ for $t = 2, \ldots, T$. This parameterization makes $q_\psi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$ in Eq. (8) tractable.

Based on Eq. (28) in Appendix, we obtain $q_\psi(\mathbf{x}_t|\mathbf{x}_{t-2})$ as follows:

$$\begin{aligned}
q_\psi(\mathbf{x}_t|\mathbf{x}_{t-2}) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2}, \big((1 - \alpha_t) + \alpha_t(1 - \alpha_{t-1})\big)\mathbf{I}) \\
&= \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2}, (1 - \alpha_t\alpha_{t-1})\mathbf{I})
\end{aligned} \tag{13}$$

By repeating the same argument, we obtain $q_\psi(\mathbf{x}_t|\mathbf{x}_0)$ as follows:

$$q_\psi(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \tag{14}$$

where $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$. It is easy to sample from $q_\psi(\mathbf{x}_t|\mathbf{x}_0)$.

We regard $\psi$ as free parameters and drop $\psi$ from our notations for the rest of this presentation. We rewrite $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ appearing in $L_{t-1}$ of Eq. (11) as follows:

$$
\begin{aligned}
q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \quad \text{(based on Eq. (8))} \\
&\propto \exp\left(-\frac{1}{2}\left(\frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{1-\bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_{t1}}\mathbf{x}_0)^2}{1-\bar{\alpha}_t}\right)\right) \\
&\propto \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right)\mathbf{x}_{t-1}^2 - 2\left(\frac{\sqrt{\alpha_t}}{1-\alpha_t}\mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}}\mathbf{x}_0\right)\mathbf{x}_{t-1}\right)\right)
\end{aligned}
\tag{15}
$$

We denote the element-wise mean and variance of $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ by $\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0)$ and $\tilde{\beta}_t$ respectively. That is,

$$
q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \equiv \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t)
\tag{16}
$$

Then we obtain

$$
\tilde{\beta}_t = 1/\left(\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right) = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{\alpha_t - \alpha_t\bar{\alpha}_{t-1} + 1 - \alpha_t} = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}(1-\alpha_t)
\tag{17}
$$

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \left( \frac{\sqrt{\alpha_t}}{1 - \alpha_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_0 \right) \Big/ \left( \frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right)$$

$$= \left( \frac{\sqrt{\alpha_t}}{1 - \alpha_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_0 \right) \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} (1 - \alpha_t)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} \mathbf{x}_0 \tag{18}$$

Based on Eq. (14), we reparameterize $\mathbf{x}_t$ as

$$\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + (1 - \bar{\alpha}_t) \boldsymbol{\epsilon} \ \text{ for } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{19}$$

and rewrite $\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0)$ as follows:

$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} \left( \frac{\mathbf{x}_t}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \boldsymbol{\epsilon} \right)$$

$$= \frac{1}{\sqrt{\alpha_t}} \left( \left( \frac{\alpha_t(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} + \frac{1 - \alpha_t}{1 - \bar{\alpha}_t} \right) \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right)$$

$$= \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right) \tag{20}$$

# Contents

# Generative modeling of observations

Here we specify the details of our Bayesian generative model firstly in this presentation.

$$p_\theta(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \tag{21}$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \tag{22}$$

We assume that $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$ as discussed in [1].

Eqs. (16) and (22) show that the KL divergence appearing in $L_{t-1}$ of Eq. (11) is from one Gaussian distribution to another. Therefore, we can rewrite $L_{t-1}$ as follows[1]:

$$L_{t-1} = -\mathbb{E}_{\neg t}\left[\frac{1}{2\sigma_t^2}\|\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2\right] + const. \tag{23}$$

---

[1] https://scoste.fr/posts/dkl_gaussian/

By using the reparameterization $\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + (1 - \bar{\alpha}_t)\boldsymbol{\epsilon}$ in Eq. (19) and the result in Eq. (20), we further rewrite $L_{t-1}$ as

$$L_{t-1} = -\mathbb{E}_{\neg t}\left[\frac{1}{2\sigma_t^2}\left\|\frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}\right) - \boldsymbol{\mu}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}), t)\right\|^2\right] + const. \qquad (24)$$

We may parameterize $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ as follows [1]:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\right) \qquad (25)$$

where $\boldsymbol{\epsilon}_\theta$ is a function approximator intended to predict $\boldsymbol{\epsilon}$ from $\mathbf{x}_t$. Then we can rewrite $L_{t-1}$ as follows:

$$L_{t-1} = -\mathbb{E}_{\neg t}\left[\frac{(1 - \alpha_t)^2}{2\sigma_t^2(1 - \bar{\alpha}_t)}\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + (1 - \bar{\alpha}_t)\boldsymbol{\epsilon}, t)\|^2\right] + const. \qquad (26)$$

We consider $L_T$ in Eq. (10):

$$L_T \equiv \int q_\psi(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_T)}{q_\psi(\mathbf{x}_T|\mathbf{x}_0)} d\mathbf{x}_{1:T} \qquad (27)$$

Both of the noise distribution $p_\theta(\mathbf{x}_T)$ and the approximate posterior $q_\psi(\mathbf{x}_T|\mathbf{x}_0)$ have no trainable parameters. Therefore, $L_T$ can be regarded as a constant.

Next, we consider $L_0$ in Eq. (10). How we maximize $L_0$ depends on how we specify the distribution $p_\theta(\mathbf{x}_0|\mathbf{x}_1)$, which directly models the observation. For example, see Sec. 3.3 of [1].

**Notice:** In this presentation, we only discuss a variational inference for diffusion models. We do not discuss where diffusion models come from.

# Contents

# Appendix

$$\int \exp\left(-\frac{(x-ay)^2}{2s^2} - \frac{(y-bz)^2}{2t^2}\right)dy = \int \exp\left(-\frac{t^2(x-ay)^2 + s^2(y-bz)^2}{2s^2t^2}\right)dy$$

$$= \int \exp\left(-\frac{(s^2+t^2a^2)y^2 - 2(s^2bz+t^2ax)y + t^2x^2 + s^2b^2z^2}{2s^2t^2}\right)dy$$

$$= \exp\left(-\frac{t^2x^2 + s^2b^2z^2}{2s^2t^2}\right)\int \exp\left(-\frac{s^2+t^2a^2}{2s^2t^2}\left(y^2 - \frac{2(s^2bz+t^2ax)}{s^2+t^2a^2}y\right)\right)dy$$

$$= \exp\left(-\frac{t^2x^2 + s^2b^2z^2}{2s^2t^2} + \frac{(s^2bz+t^2ax)^2}{2s^2t^2(s^2+t^2a^2)}\right)\int \exp\left(-\frac{s^2+t^2a^2}{2s^2t^2}\left(y - \frac{s^2bz+t^2ax}{s^2+t^2a^2}\right)^2\right)dy$$

$$\propto \exp\left(-\frac{s^2t^2x^2 + s^4b^2z^2 + t^4a^2x^2 + s^2t^2a^2b^2z^2 - t^4a^2x^2 - 2s^2t^2abzx - s^4b^2z^2}{2s^2t^2(s^2+t^2a^2)}\right)$$

$$= \exp\left(-\frac{x^2 - 2abzx + a^2b^2z^2}{2(s^2+t^2a^2)}\right) = \exp\left(-\frac{(x-abz)^2}{2(s^2+t^2a^2)}\right) \tag{28}$$

📄 Jonathan Ho, Ajay Jain, and Pieter Abbeel.

Denoising diffusion probabilistic models.

*CoRR*, abs/2006.11239, 2020.