

A Collection of Frequently Asked Questions about Denoising Diffusion Probabilistic Models*

Kaibo Tang

November 21, 2024

1 Does the forward process tend to $\mathcal{N}(\mathbf{0}, \mathbf{I})$?

Proposition 1.1. *Suppose $a_n \in (0, 1) \forall n \in \mathbb{N}$. Then,*

$$\prod_{n=1}^{\infty} (1 - a_n) = 0 \iff \sum_{n=1}^{\infty} a_n \text{ diverges.} \quad (1)$$

Proof.

1. (\implies). For the sake of contradiction, assume $\sum_{n=1}^{\infty} a_n$ converges, in which case, we know $a_n \rightarrow 0$ as $n \rightarrow \infty$. In addition, since

$$\prod_{n=1}^{\infty} (1 - a_n) = 0, \quad (2)$$

we have

$$\sum_{n=1}^{\infty} -\ln(1 - a_n) = \infty. \quad (3)$$

By L'Hopital's rule, we see that

$$\lim_{n \rightarrow \infty} -\frac{a_n}{\ln(1 - a_n)} = \lim_{n \rightarrow \infty} (1 - a_n) = 1, \quad (4)$$

and limit comparison test applies, in which case, we know $\sum_{n=1}^{\infty} a_n$ also diverges, which is contradiction.

*<https://arxiv.org/abs/2006.11239>

2. (\Leftarrow). Notice that

$$\sum_{n=1}^{\infty} -\ln(1 - a_n) \geq \sum_{n=1}^{\infty} a_n. \quad (5)$$

By comparison test, we see that

$$\sum_{n=1}^{\infty} -\ln(1 - a_n) = \infty, \quad (6)$$

i.e.,

$$\prod_{n=1}^{\infty} (1 - a_n) = 0. \quad (7)$$

□

Proposition 1.2. *Given $\beta_t \in (0, 1)$ such that $\sum_{t=1}^{\infty} \beta_t$ diverges, a Markov chain with Gaussian transitions*

$$p(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (8)$$

satisfies

$$p(\mathbf{x}_t) \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \text{as } t \rightarrow \infty. \quad (9)$$

Proof.

1. For each $t \geq 1$, let

$$\mathbf{y}_t = \frac{\mathbf{x}_t - \sqrt{1 - \beta_t} \mathbf{x}_{t-1}}{\sqrt{\beta_t}}. \quad (10)$$

We see that \mathbf{y}_t is a standard normal random variable, given $\mathbf{x}_0, \dots, \mathbf{x}_{t-1}$. Hence, we see that \mathbf{y}_t is independent on $\mathbf{x}_0, \dots, \mathbf{x}_{t-1}$. Similarly, we see that $\mathbf{x}_0, \mathbf{y}_1, \dots, \mathbf{y}_t$ are independent.

2. Now, notice that

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathbf{y}_t \quad (11)$$

$$= \sqrt{1 - \beta_t} \sqrt{1 - \beta_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \beta_t} \sqrt{\beta_{t-1}} \mathbf{y}_{t-1} + \sqrt{\beta_t} \mathbf{y}_t \quad (12)$$

$$= \dots \quad (13)$$

$$= \prod_{i=1}^t \sqrt{1 - \beta_i} \mathbf{x}_0 + \sum_{i=1}^{t-1} \left[\left(\prod_{j=i+1}^t \sqrt{1 - \beta_j} \right) \sqrt{\beta_i} \mathbf{y}_i \right] + \sqrt{\beta_t} \mathbf{y}_t. \quad (14)$$

We see that the distribution of \mathbf{x}_t given \mathbf{x}_0 is normal with expectation

$$\boldsymbol{\mu} = \prod_{i=1}^t \sqrt{1 - \beta_i} \mathbf{x}_0, \quad (15)$$

and covariance matrix

$$\boldsymbol{\Sigma} = \sum_{i=1}^{t-1} \left[\left(\prod_{j=i+1}^t (1 - \beta_j) \right) \beta_i \mathbf{I} \right] + \beta_t \mathbf{I}. \quad (16)$$

To simplify the covariance matrix, notice that

$$\left(\prod_{j=i+1}^t (1 - \beta_j) \right) \beta_i = \left(\prod_{j=i+1}^t (1 - \beta_j) \right) (1 - (1 - \beta_i)) \quad (17)$$

$$= \prod_{j=i+1}^t (1 - \beta_j) - \prod_{j=i}^t (1 - \beta_j). \quad (18)$$

Hence, we have

$$\boldsymbol{\Sigma} = \sum_{i=1}^{t-1} \left[\prod_{j=i+1}^t (1 - \beta_j) \mathbf{I} - \prod_{j=i}^t (1 - \beta_j) \mathbf{I} \right] + \beta_t \mathbf{I} \quad (19)$$

$$= (1 - \beta_t) \mathbf{I} - \prod_{j=1}^t (1 - \beta_j) \mathbf{I} + \beta_t \mathbf{I} \quad (20)$$

$$= \left[1 - \prod_{j=1}^t (1 - \beta_j) \right] \mathbf{I}. \quad (21)$$

3. Finally, since $\sum_{j=1}^{\infty} \beta_j$ diverges, we see that $\prod_{j=1}^{\infty} (1 - \beta_j) = 0$ from Proposition 1.1. Hence, we see that as $t \rightarrow \infty$, we have

$$\boldsymbol{\mu} = \prod_{i=1}^t \sqrt{1 - \beta_i} \mathbf{x}_0 \rightarrow \mathbf{0}, \quad \boldsymbol{\Sigma} = \left[1 - \prod_{j=1}^t (1 - \beta_j) \right] \mathbf{I} \rightarrow \mathbf{I} \quad (22)$$

for any \mathbf{x}_0 finite.

□

Remark 1.1. Let $\alpha_j := 1 - \beta_j$ and $\bar{\alpha}_t := \prod_{j=1}^t \alpha_j$. We see that

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}). \quad (23)$$

2 How do you get from NLL to L_{simple} ?

Starting with the negative log likelihood,

$$-\ln p_{\theta}(\mathbf{x}_0), \quad (24)$$

the authors manage to reduce this objective to

$$L_{\text{simple}}(\theta) = E_{t,\epsilon} \left(\left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t) \right\|^2 \right). \quad (25)$$

This is done in three steps, which I will outline below.

2.1 ELBO

We first find a variational bound for the negative log likelihood term, i.e.,

$$-\ln p_{\theta}(\mathbf{x}_0) \leq E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left(-\ln \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right). \quad (26)$$

Proof of this inequality involves noticing that $\ln(\cdot)$ is a convex function and invoking Jensen's inequality.

2.2 KL divergence

Using the definition of p_{θ} and q , we see that

$$-\ln \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} = -\ln p(\mathbf{x}_T) - \sum_{t=1}^T \ln \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \quad (27)$$

$$= -\ln p(\mathbf{x}_T) - \sum_{t=2}^T \ln \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} - \ln \frac{p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)}, \quad (28)$$

where the second equality follows from taking the first term out of the summation. Notice that

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})q(\mathbf{x}_{t-1}|\mathbf{x}_0) = q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0). \quad (29)$$

Hence, we have

$$-\ln \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} = -\ln p(\mathbf{x}_T) - \sum_{t=2}^T \ln \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \cdot \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} - \ln \frac{p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} \quad (30)$$

$$= -\ln \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} - \sum_{t=2}^T \ln \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} - \ln p_\theta(\mathbf{x}_0|\mathbf{x}_1). \quad (31)$$

Now, we notice that

$$L_T := E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left(-\ln \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right) \quad (32)$$

$$= E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[E_{\mathbf{x}_T \sim q(\mathbf{x}_T|\mathbf{x}_0)} \left(-\ln \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right) \right] \quad (33)$$

$$= E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\int q(\mathbf{x}_T|\mathbf{x}_0) \ln \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p(\mathbf{x}_T)} d\mathbf{x}_T \right] \quad (34)$$

$$= E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))]. \quad (35)$$

Remark 2.1. Given non-learnable variance schedule β_t , the term L_T has no learnable parameter and is thus constant during training.

Notice that

$$E_{\mathbf{x}_{t-1}, \mathbf{x}_t \sim q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0)} \left(-\ln \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right) \quad (36)$$

$$= \iint q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0) \ln \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} d\mathbf{x}_{t-1} d\mathbf{x}_t \quad (37)$$

$$= \int q(\mathbf{x}_t|\mathbf{x}_0) \int q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \ln \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} d\mathbf{x}_{t-1} d\mathbf{x}_t \quad (38)$$

$$= E_{\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))], \quad (39)$$

where we use the fact that

$$q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0) \quad (40)$$

for the second equality. Hence, we have

$$L_{t-1} := E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left(-\ln \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right) \quad (41)$$

$$= E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[E_{\mathbf{x}_{t-1}, \mathbf{x}_t \sim q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0)} \left(-\ln \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right) \right] \quad (42)$$

$$= E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] \quad (43)$$

Finally, let

$$L_0 := E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} (-\ln p_\theta(\mathbf{x}_0|\mathbf{x}_1)). \quad (44)$$

We see that now the variational bound is given by

$$-\ln p_\theta(\mathbf{x}_0) \leq L_T + \sum_{t=2}^T L_{t-1} + L_0. \quad (45)$$

2.3 Further simplification

We state (without proof) the expression for the KL divergence between two multivariate normal distributions.¹

Lemma 2.1. *Let P and Q denote two multivariate normal distributions with means $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^n$, respectively, and covariance matrices $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathcal{M}_{n \times n}$, respectively. Then, the KL divergence of P from Q is given by*

$$D_{\text{KL}}(P \parallel Q) = \frac{1}{2} \left[(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) - \ln \frac{\det \boldsymbol{\Sigma}_1}{\det \boldsymbol{\Sigma}_2} - n \right]. \quad (46)$$

Corollary 2.1. *In the previous case, let $\boldsymbol{\Sigma}_1 = \sigma_1 \mathbf{I}$ and $\boldsymbol{\Sigma}_2 = \sigma_2 \mathbf{I}$. Then,*

$$D_{\text{KL}}(P \parallel Q) = \frac{1}{2} \left[\frac{\|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|^2}{\sigma_2^2} + \frac{n\sigma_1^2}{\sigma_2^2} - \ln \frac{\sigma_1}{\sigma_2} - n \right]. \quad (47)$$

Now, consider the reverse process at time t given by

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \quad \text{for } t = 2, 3, \dots, T. \quad (48)$$

¹A detailed proof can be found at <https://statproofbook.github.io/P/mvn-kl.html>

The following simplification assumes the covariance matrix at each step to be non-learnable and dependent only on time t , i.e.,

$$\Sigma_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}, \quad (49)$$

in which case, we have

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}) \quad (50)$$

Recall that for the forward process posterior, we have

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}), \quad (51)$$

$$\text{where } \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \quad (52)$$

$$\text{and } \tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t. \quad (53)$$

From Corollary 2.1, we see that

$$D_{\text{KL}}(P \parallel Q) = \frac{1}{2\sigma_t^2} \|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2 + \text{const.}, \quad (54)$$

since σ_t and $\tilde{\beta}_t$ are fixed.

Now, we recall equation (23), i.e.,

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}). \quad (55)$$

Notice that we can reparameterize \mathbf{x}_t as

$$\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (56)$$

Substituting (52), (56) into (54), we see that

$$L_{t-1} = E_\epsilon \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right) - \boldsymbol{\mu}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}), t) \right\|^2 \right] + \text{const.} \quad (57)$$

Instead of directly parameterizing $\boldsymbol{\mu}_\theta$, consider

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right), \quad (58)$$

where $\boldsymbol{\epsilon}_\theta$ predicts $\boldsymbol{\epsilon}$ from \mathbf{x}_t .

Remark 2.2. During inference, we need to sample $\mathbf{x}_{t-1} \sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$. Reparameterizing again, we see that

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z} \quad \text{where } \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (59)$$

Finally, we arrive at

$$L_{t-1} = E_\epsilon \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2 \right] \quad (60)$$

3 Disclaimer

There are still several loose ends that are yet to be tied. For example, we have not demonstrated how to deal with L_0 , which is nontrivial. But I believe the most important stuff are well-covered in this note. In addition, I would like to also use this opportunity to highlight some caveats.

1. Although the authors do not make this explicit, in Proposition 1.2, we see that the divergence of $\sum_{t=1}^{\infty} \beta_t$ is essential for DDPM to work.
2. In the original paper, the authors start with $E_{\mathbf{x}_0}(-\ln p_\theta(\mathbf{x}_0))$. However, I think taking the expectation over \mathbf{x}_0 creates confusion. Hence, I started with $-\ln p_\theta(\mathbf{x}_0)$.
3. The authors use E_q ambiguously. Since I dropped $E_{\mathbf{x}_0}$, we only need to deal with $E_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)}$ and later on E_ϵ , which makes the derivation a lot clearer.
4. Proposition 1.1, Proposition 1.2, Lemma 2.1, and Corollary 2.1 are great conclusions to remember.