

# Assignment 1: Normalizing Flows & Diffusion Models

Deep Generative Models Course 361.2.2370

December 10, 2024

## Contents

<b>1</b>	<b>Diffusion Models</b>	<b>1</b>
1.1	Conditional Diffusion Distribution . . . . .	1
1.1.1	(a) . . . . .	1
1.1.2	(b) . . . . .	3
1.1.3	(c) . . . . .	4
1.1.4	(d) . . . . .	4
1.2	Evidence Lower Bound (ELBO) . . . . .	6
1.2.1	(a) . . . . .	6
1.2.2	(b) . . . . .	8
1.2.3	(c) . . . . .	10

## 1 Diffusion Models

In this exercise, we will drill down into the probabilistic framework underlying denoising diffusion models.

### 1.1 Conditional Diffusion Distribution

We defined the conditional probability  $q(\mathbf{z}_t|\mathbf{z}_{t-1})$  as the mixing process. To reverse this process, we apply Bayes' rule:

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t) = \frac{q(\mathbf{z}_t|\mathbf{z}_{t-1})q(\mathbf{z}_{t-1})}{q(\mathbf{z}_t)}$$

This is intractable since we cannot compute the marginal distribution  $q(\mathbf{z}_{t-1})$ . However, if we know the starting variable  $\mathbf{x}$ , we do know the distribution  $q(\mathbf{z}_{t-1}|\mathbf{x})$ .

#### 1.1.1 (a)

Consider two multivariate statistically independent normal distributions,

$$P_1(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \quad P_2(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2).$$

## Answer

The multivariate Gaussian distribution is given by:

$$P(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right).$$

The product of  $P_1(\mathbf{x})$  and  $P_2(\mathbf{x})$  is:

$$P_1(\mathbf{x}) \cdot P_2(\mathbf{x}) = \frac{1}{(2\pi)^d |\boldsymbol{\Sigma}_1|^{1/2} |\boldsymbol{\Sigma}_2|^{1/2}} \exp \left( -\frac{1}{2} \left[ (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + (\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right] \right).$$

Expanding the quadratic terms:

$$(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) = \mathbf{x}^T \boldsymbol{\Sigma}_1^{-1} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1,$$

$$(\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) = \mathbf{x}^T \boldsymbol{\Sigma}_2^{-1} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2.$$

Combining the terms:

$$P_1(\mathbf{x}) \cdot P_2(\mathbf{x}) \propto \exp \left( -\frac{1}{2} \left[ \mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1}) \mathbf{x} - 2\mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) \right] \right).$$

We will compare the resulting expression to the expression for a normal distribution:

$$\mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1}) \mathbf{x} - 2\mathbf{x}^T (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) = \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$$

And we get:

$$\begin{aligned} \boldsymbol{\Sigma}_{\text{new}}^{-1} &= \boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1}, \\ \boldsymbol{\Sigma}_{\text{new}}^{-1} \boldsymbol{\mu}_{\text{new}} &= \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2. \end{aligned}$$

Multiplying  $\boldsymbol{\Sigma}_{\text{new}}$  to get  $\boldsymbol{\mu}_{\text{new}}$ :

$$\boldsymbol{\mu}_{\text{new}} = \boldsymbol{\Sigma}_{\text{new}} (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2).$$

Therefore, the covariance and mean of the resulting distribution are:

$$\begin{aligned} \boldsymbol{\Sigma}_{\text{new}} &= (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}, \\ \boldsymbol{\mu}_{\text{new}} &= (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1} (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2). \end{aligned}$$

Thus, the product is proportional to:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\text{new}}, \boldsymbol{\Sigma}_{\text{new}}).$$

Substituting these results yields the desired form:

$$P_1(\mathbf{x}) \cdot P_2(\mathbf{x}) \propto \mathcal{N}(\mathbf{x}; (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1} (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2), (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}).$$

Proves the relation.

### 1.1.2 (b)

Consider a multivariate normal distribution  $\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B})$ . Prove the following change of variable identity:

$$\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B}) \propto \mathcal{N}(\mathbf{y}; (A^T \mathbf{B}^{-1} A)^{-1} A^T \mathbf{B}^{-1} \mathbf{x}, (A^T \mathbf{B}^{-1} A)^{-1}).$$

### Proof

The probability density function of the multivariate Gaussian  $\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B})$  is given by:

$$\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B}) = \frac{1}{(2\pi)^{d/2} |\mathbf{B}|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - A \cdot \mathbf{y})^T \mathbf{B}^{-1} (\mathbf{x} - A \cdot \mathbf{y}) \right).$$

Expanding the quadratic term in the exponent:

$$(\mathbf{x} - A \cdot \mathbf{y})^T \mathbf{B}^{-1} (\mathbf{x} - A \cdot \mathbf{y}) = \mathbf{x}^T \mathbf{B}^{-1} \mathbf{x} - 2\mathbf{x}^T \mathbf{B}^{-1} A \cdot \mathbf{y} + \mathbf{y}^T A^T \mathbf{B}^{-1} A \cdot \mathbf{y}.$$

The first term,  $\mathbf{x}^T \mathbf{B}^{-1} \mathbf{x}$ , is a constant with respect to  $\mathbf{y}$ , so we can write:

$$\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B}) \propto \exp \left( -\frac{1}{2} (-2\mathbf{x}^T \mathbf{B}^{-1} A \cdot \mathbf{y} + \mathbf{y}^T A^T \mathbf{B}^{-1} A \cdot \mathbf{y}) \right).$$

Rewriting the exponent:

$$\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B}) \propto \exp \left( -\frac{1}{2} (\mathbf{y}^T A^T \mathbf{B}^{-1} A \cdot \mathbf{y} - 2\mathbf{y}^T A^T \mathbf{B}^{-1} \mathbf{x}) \right).$$

We will compare the resulting expression

$$\mathbf{y}^T A^T \mathbf{B}^{-1} A \cdot \mathbf{y} - 2\mathbf{y}^T A^T \mathbf{B}^{-1} \mathbf{x},$$

to the expression for a normal distribution:

$$\mathbf{y}^T \Sigma^{-1} \mathbf{y} - 2\mathbf{y}^T \Sigma^{-1} \boldsymbol{\mu}$$

And we get:

$$\begin{aligned} \Sigma_y^{-1} &= A^T \mathbf{B}^{-1} A, \\ \Sigma_y^{-1} \boldsymbol{\mu}_y &= A^T \mathbf{B}^{-1} \mathbf{x}. \end{aligned}$$

Multiplying  $\Sigma_y$  to get  $\boldsymbol{\mu}_y$ :

$$\boldsymbol{\mu}_y = \Sigma_y (A^T \mathbf{B}^{-1} \mathbf{x}).$$

Therefore, the covariance and mean of the resulting distribution are:

$$\begin{aligned} \boldsymbol{\mu}_y &= (A^T \mathbf{B}^{-1} A)^{-1} A^T \mathbf{B}^{-1} \mathbf{x}, \\ \Sigma_y &= (A^T \mathbf{B}^{-1} A)^{-1}. \end{aligned}$$

Thus, we have shown:

$$\mathcal{N}(\mathbf{x}; A \cdot \mathbf{y}, \mathbf{B}) \propto \mathcal{N}(\mathbf{y}; (A^T \mathbf{B}^{-1} A)^{-1} A^T \mathbf{B}^{-1} \mathbf{x}, (A^T \mathbf{B}^{-1} A)^{-1}).$$

### 1.1.3 (c)

Show that:

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) \propto q(\mathbf{z}_t | \mathbf{z}_{t-1}) \cdot q(\mathbf{z}_{t-1} | \mathbf{x}).$$

### Answer

Using Bayes rule:

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) = \frac{q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}) \cdot q(\mathbf{z}_{t-1} | \mathbf{x})}{q(\mathbf{z}_t | \mathbf{x})}.$$

Assuming conditional independence of  $\mathbf{z}_t$  and  $\mathbf{x}$  given  $\mathbf{z}_{t-1}$ :

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

We substitute this into the equation:

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) = \frac{q(\mathbf{z}_t | \mathbf{z}_{t-1}) \cdot q(\mathbf{z}_{t-1} | \mathbf{x})}{q(\mathbf{z}_t | \mathbf{x})}.$$

Here,  $q(\mathbf{z}_t | \mathbf{x})$  integrates over all possible values of  $\mathbf{z}_{t-1}$ , thus becoming a normalization constant that does not depend on  $\mathbf{z}_{t-1}$ . Therefore, we can express the proportionality:

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) \propto q(\mathbf{z}_t | \mathbf{z}_{t-1}) \cdot q(\mathbf{z}_{t-1} | \mathbf{x}).$$

### 1.1.4 (d)

Use (a), (b), and (c) to find a proportional term for  $q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})$ , where:

- $q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \cdot \mathbf{z}_{t-1}, \beta_t \cdot \mathbf{I}),$
- $q(\mathbf{z}_{t-1} | \mathbf{x}) = \mathcal{N}(\mathbf{z}_{t-1}; \sqrt{\alpha_{t-1}} \cdot \mathbf{x}, (1 - \alpha_{t-1}) \cdot \mathbf{I}).$

**Hint:** Note that the relation in (a) requires the normal distributions to have the same support.

### Answer

From (c) we know:

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) \propto q(\mathbf{z}_t | \mathbf{z}_{t-1}) \cdot q(\mathbf{z}_{t-1} | \mathbf{x})$$

and we know that  $q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \cdot \mathbf{z}_{t-1}, \beta_t \cdot \mathbf{I}).$

we will denote  $A = \sqrt{1 - \beta_t}$ ,  $B = \beta_t \cdot \mathbf{I}$ ,  $\mathbf{x} = \mathbf{z}_t$ ,  $\mathbf{y} = \mathbf{z}_{t-1}$ , and we will calculate the expressions from (b):

$$\mu_{\mathbf{z}_{t-1}} = (A^T \mathbf{B}^{-1} A)^{-1} A^T \mathbf{B}^{-1} \mathbf{z}_t = \frac{\sqrt{1 - \beta_t} (\beta_t \cdot \mathbf{I})^{-1}}{\sqrt{1 - \beta_t} (\beta_t \cdot \mathbf{I})^{-1} \sqrt{1 - \beta_t}} \mathbf{z}_t = \frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t$$

$$\Sigma_{z_{-1}} = (A^T B^{-1} A)^{-1} = \frac{1}{\sqrt{1 - \beta_t}(\beta_t \cdot I)^{-1} \sqrt{1 - \beta_t}} = \frac{\beta_t}{1 - \beta_t} \cdot I$$

Therefore:

$$q(z_{t-1} \mid z_t) \propto \mathcal{N}\left(z_{t-1}; \frac{1}{\sqrt{1 - \beta_t}} z_t, \frac{\beta_t}{1 - \beta_t} I\right)$$

We also know:

$$q(z_{t-1} \mid x) = \mathcal{N}(z_{t-1}; \sqrt{\alpha_{t-1}} \cdot x, (1 - \alpha_{t-1}) \cdot I)$$

Combining  $q(z_{t-1} \mid z_t)$  and  $q(z_{t-1} \mid x)$ , both expressed as Gaussian distributions in  $z_{t-1}$ , we have:

$$q(z_{t-1} \mid z_t, x) \propto \mathcal{N}(z_{t-1}; \mu_1, \Sigma_1) \cdot \mathcal{N}(z_{t-1}; \mu_2, \Sigma_2)$$

where:

$$\begin{aligned} \mu_1 &= \frac{1}{\sqrt{1 - \beta_t}} z_t, & \Sigma_1 &= \frac{\beta_t}{1 - \beta_t} I \\ \mu_2 &= \sqrt{\alpha_{t-1}} \cdot x, & \Sigma_2 &= (1 - \alpha_{t-1}) \cdot I \end{aligned}$$

From (a) we know that the covariance of the resulting Gaussian is:

$$\Sigma = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}$$

Substitute:

$$\begin{aligned} \Sigma_1^{-1} &= \frac{1 - \beta_t}{\beta_t} I, & \Sigma_2^{-1} &= \frac{1}{1 - \alpha_{t-1}} I \\ \Sigma &= \left( \frac{1 - \beta_t}{\beta_t} + \frac{1}{1 - \alpha_{t-1}} \right)^{-1} I \end{aligned}$$

Simplify:

$$\begin{aligned} \frac{1 - \beta_t}{\beta_t} + \frac{1}{1 - \alpha_{t-1}} &= \frac{(1 - \beta_t)(1 - \alpha_{t-1}) + \beta_t}{\beta_t(1 - \alpha_{t-1})} \\ &= \frac{1 - \alpha_{t-1} - \beta_t - \beta_t \alpha_{t-1} + \beta_t}{\beta_t(1 - \alpha_{t-1})} \\ &= \frac{1 - \alpha_{t-1}(1 - \beta_t)}{\beta_t(1 - \alpha_{t-1})} \\ &= \frac{1 - \alpha_t}{\beta_t(1 - \alpha_{t-1})} \end{aligned}$$

Where The last equality follows from the definition of  $\alpha_t$ :

$$\alpha_t = \alpha_{t-1}(1 - \beta_t)$$

Thus:

$$\Sigma = \left( \frac{1 - \alpha_t}{\beta_t(1 - \alpha_{t-1})} \right)^{-1} = \frac{\beta_t(1 - \alpha_{t-1})}{1 - \alpha_t} I$$

The mean of the resulting Gaussian is:

$$\mu = \Sigma (\Sigma_1^{-1} \mu_1 + \Sigma_2^{-1} \mu_2)$$

Substitute:

$$\begin{aligned}\Sigma_1^{-1} \mu_1 &= \frac{1 - \beta_t}{\beta_t} \cdot \frac{1}{\sqrt{1 - \beta_t}} z_t = \frac{\sqrt{1 - \beta_t}}{\beta_t} z_t \\ \Sigma_2^{-1} \mu_2 &= \frac{1}{1 - \alpha_{t-1}} \cdot \sqrt{\alpha_{t-1}} x = \frac{\sqrt{\alpha_{t-1}}}{1 - \alpha_{t-1}} x \\ \mu &= \frac{\beta_t(1 - \alpha_{t-1})}{1 - \alpha_t} \left( \frac{\sqrt{1 - \beta_t}}{\beta_t} z_t + \frac{\sqrt{\alpha_{t-1}}}{1 - \alpha_{t-1}} x \right)\end{aligned}$$

Simplify:

$$\mu = \frac{(1 - \alpha_{t-1})\sqrt{1 - \beta_t}}{1 - \alpha_t} z_t + \frac{\beta_t\sqrt{\alpha_{t-1}}}{1 - \alpha_t} x$$

## Final Expression

$$q(z_{t-1} \mid z_t, x) = \mathcal{N}(z_{t-1}; \mu, \Sigma)$$

Where:

$$\begin{aligned}\mu &= \frac{(1 - \alpha_{t-1})\sqrt{1 - \beta_t}}{1 - \alpha_t} z_t + \frac{\beta_t\sqrt{\alpha_{t-1}}}{1 - \alpha_t} x \\ \Sigma &= \frac{\beta_t(1 - \alpha_{t-1})}{1 - \alpha_t} I\end{aligned}$$

## 1.2 Evidence Lower Bound (ELBO)

### 1.2.1 (a)

Consider two  $D$ -dimensional multivariate statistically independent normal distributions,  $N(x; \mu_1, \Sigma_1)$  and  $N(x; \mu_2, \Sigma_2)$ . Derive a concise closed-form formula for the Kullback-Leibler (KL) divergence:

$$D_{KL}(N(x; \mu_1, \Sigma_1) \parallel N(x; \mu_2, \Sigma_2)).$$

## Answer

Let

$$P := x \sim \mathcal{N}(\mu_1, \Sigma_1), \quad Q := x \sim \mathcal{N}(\mu_2, \Sigma_2)$$

The KL divergence for a continuous random variable is defined as:

$$D_{KL}[P \parallel Q] = \int_X p(x) \log \frac{p(x)}{q(x)} dx.$$

Substituting the probability density functions of the multivariate normal distributions, we get:

$$\begin{aligned} D_{KL}[P \parallel Q] &= \int_{\mathbb{R}^n} \mathcal{N}(x; \mu_1, \Sigma_1) \log \frac{\mathcal{N}(x; \mu_1, \Sigma_1)}{\mathcal{N}(x; \mu_2, \Sigma_2)} dx \\ &= \mathbb{E}_P \left[ \log \frac{\mathcal{N}(x; \mu_1, \Sigma_1)}{\mathcal{N}(x; \mu_2, \Sigma_2)} \right] \end{aligned}$$

The log term can be expanded as:

$$\begin{aligned} \log \frac{\mathcal{N}(x; \mu_1, \Sigma_1)}{\mathcal{N}(x; \mu_2, \Sigma_2)} &= \log \frac{\frac{1}{\sqrt{(2\pi)^n |\Sigma_1|}}}{\frac{1}{\sqrt{(2\pi)^n |\Sigma_2|}}} - \frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \\ &= \log \frac{\sqrt{(2\pi)^n}}{\sqrt{(2\pi)^n}} + \log \sqrt{\frac{|\Sigma_2|}{|\Sigma_1|}} - \frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \\ &= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \end{aligned}$$

Simplifying:

$$\log \frac{\mathcal{N}(x; \mu_1, \Sigma_1)}{\mathcal{N}(x; \mu_2, \Sigma_2)} = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right].$$

We will use the following properties:

1. Since the terms  $(x - \mu) \Sigma^{-1} (x - \mu)$  are scalars,  $(x - \mu) \Sigma^{-1} (x - \mu) = \mathbf{tr}[(x - \mu) \Sigma^{-1} (x - \mu)]$ .
2.  $\mathbf{tr}[ABC] = \mathbf{tr}[BCA]$ .
3.  $\mathbb{E}[X^T X] = \mu^T \mu + \mathbf{tr}(\Sigma)$ .

And we get:

$$\begin{aligned}
D_{KL}[P \parallel Q] &= \mathbb{E}_P \left[ \log \frac{\mathcal{N}(x; \mu_1, \Sigma_1)}{\mathcal{N}(x; \mu_2, \Sigma_2)} \right] = \\
&= \mathbb{E}_P \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right] \\
&= \mathbb{E}_P \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - \mathbf{tr} [(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)] + \mathbf{tr} [(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)] \right] \\
&= \mathbb{E}_P \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - \mathbf{tr} [\Sigma_1^{-1} (x - \mu_1)^T (x - \mu_1)] + \mathbf{tr} [\Sigma_2^{-1} (x - \mu_2)^T (x - \mu_2)] \right] \\
&= \mathbb{E}_P \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} + \mathbf{tr} [\Sigma_1^{-1} [-xx^T + 2x\mu_1 - \mu_1^T \mu_1]] + \mathbf{tr} [\Sigma_2^{-1} [xx^T - 2x\mu_2 + \mu_2^T \mu_2]] \right] \\
&= \mathbb{E}_P \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} + \mathbf{tr} [(-\Sigma_1^{-1} + \Sigma_2^{-1})xx^T + 2x[\Sigma_1^{-1}\mu_1 - \Sigma_2^{-1}\mu_2]] - \Sigma_1^{-1}\mu_1^T \mu_1 + \Sigma_2^{-1}\mu_2^T \mu_2 \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} + \mathbf{tr} [-\Sigma_1^{-1} + \Sigma_2^{-1}] \mathbb{E}_P [xx^T] + \mathbb{E}_P [2x] \mathbf{tr} [\Sigma_1^{-1}\mu_1 - \Sigma_2^{-1}\mu_2] - \mathbf{tr} [\Sigma_1^{-1}\mu_1^T \mu_1 + \Sigma_2^{-1}\mu_2^T \mu_2] \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} + \mathbf{tr} [-\Sigma_1^{-1} + \Sigma_2^{-1}] [\mu_1 \mu_1^T + \mathbf{tr}(\Sigma_1)] + 2\mu_1 \mathbf{tr} [\Sigma_1^{-1}\mu_1 - \Sigma_2^{-1}\mu_2] - \mathbf{tr} [\Sigma_1^{-1}\mu_1^T \mu_1 + \Sigma_2^{-1}\mu_2^T \mu_2] \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} + \mathbf{tr} [-\Sigma_1^{-1}\Sigma_1 + \Sigma_2^{-1}\Sigma_1] + [-\Sigma_1^{-1} + \Sigma_2^{-1}] \mu_1 \mu_1^T + 2\mu_1 \mathbf{tr} [\Sigma_1^{-1}\mu_1 - \Sigma_2^{-1}\mu_2] - \mathbf{tr} [\Sigma_1^{-1}\mu_1^T \mu_1 + \Sigma_2^{-1}\mu_2^T \mu_2] \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \mathbf{tr} [\Sigma_2^{-1}\Sigma_1 - \Sigma_1^{-1}\mu_1 \mu_1^T + \Sigma_2^{-1}\mu_1 \mu_1^T + \Sigma_1^{-1}2\mu_1 \mu_1^T - \Sigma_2^{-1}2\mu_1 \mu_2 - \Sigma_1^{-1}\mu_1^T \mu_1 + \Sigma_2^{-1}\mu_2^T \mu_2] \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \mathbf{tr}(\Sigma_2^{-1}\Sigma_1) + \mathbf{tr} [\Sigma_2^{-1}\mu_1 \mu_1^T - \Sigma_2^{-1}2\mu_1 \mu_2^T + \Sigma_2^{-1}\mu_2 \mu_2^T] \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \mathbf{tr}(\Sigma_2^{-1}\Sigma_1) + \mathbf{tr} [\mu_1^T \Sigma_2^{-1} \mu_1 - 2\mu_2^T \Sigma_2^{-1} \mu_1 + \mu_2^T \Sigma_2^{-1} \mu_2] \right] \\
&= \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \mathbf{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2 (\mu_2 - \mu_1) \right]
\end{aligned}$$

Final Expression:

$$D_{KL}[P \parallel Q] = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - n + \mathbf{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right].$$

### 1.2.2 (b)

Use (a) to derive an expression for:

$$D_{KL}(q(z_{t-1} \mid z_t, x) \parallel P(z_{t-1} \mid z_t; \phi_t)),$$

where:

- $q(z_{t-1} \mid z_t, x) = N \left( z_{t-1}; \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot \sqrt{1-\beta_t} \cdot z_t + \frac{\beta_t \sqrt{\alpha_{t-1}}}{1-\alpha_t} \cdot x, \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot I \right)$
- $P(z_{t-1} \mid z_t; \phi_t) = N(z_{t-1}; f_{\phi_t}(z_t), \sigma_t^2 \cdot I)$



## Answer

The goal is to derive an expression for:

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

Where:

- $\mu_1 = \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot \sqrt{1-\beta_t} \cdot z_t + \frac{\beta_t \sqrt{\alpha_{t-1}}}{1-\alpha_t} \cdot x$
- $\mu_2 = f_{\phi_t}(z_t)$
- $\Sigma_1 = \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot I$
- $\Sigma_2 = \sigma_t^2 \cdot I$

Based on (a), the KL divergence between two multivariate Gaussian distributions  $q \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $p \sim \mathcal{N}(\mu_2, \Sigma_2)$  is given by:

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

where  $n$  is the dimensionality of the distributions. The covariance matrices are diagonal (scaled identity matrices), so their determinants are the product of their diagonal entries:

$$|\Sigma_1| = \left( \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \right)^n, \quad |\Sigma_2| = (\sigma_t^2)^n$$

Thus, the log term is:

$$\begin{aligned} \log \frac{|\Sigma_2|}{|\Sigma_1|} &= \log(\sigma_t^2)^n - \log \left( \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \right)^n \\ &= n \log \sigma_t^2 - n \log \left( \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \right) \end{aligned}$$

To derive the trace term:

$$\text{tr}(\Sigma_2^{-1}\Sigma_1),$$

Since both  $\Sigma_1$  and  $\Sigma_2$  are diagonal:

$$\begin{aligned} \Sigma_2^{-1} &= \frac{1}{\sigma_t^2} I \\ \text{tr}(\Sigma_2^{-1}\Sigma_1) &= \text{tr} \left( \frac{1}{\sigma_t^2} \cdot \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot I \right) \\ &= \frac{\beta_t \cdot (1-\alpha_{t-1})}{\sigma_t^2(1-\alpha_t)} \cdot n \end{aligned}$$

The quadratic term is:

$$(\mu_1 - \mu_2)^T \Sigma_2^{-1}(\mu_1 - \mu_2)$$

Substitute  $\mu_1 = \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot \sqrt{1-\beta_t} \cdot z_t + \frac{\beta_t\sqrt{\alpha_{t-1}}}{1-\alpha_t} \cdot x$  and  $\mu_2 = f_{\phi_t}(z_t)$ :

$$\mu_1 - \mu_2 = \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot \sqrt{1-\beta_t} \cdot z_t + \frac{\beta_t\sqrt{\alpha_{t-1}}}{1-\alpha_t} \cdot x - f_{\phi_t}(z_t)$$

Let  $\delta = \mu_1 - \mu_2$ . Then:

$$\delta^T \Sigma_2^{-1} \delta = \delta^T \cdot \frac{1}{\sigma_t^2} \cdot \delta = \frac{1}{\sigma_t^2} \delta^T \cdot \delta = \frac{1}{\sigma_t^2} \cdot \|\delta\|^2$$

Where  $\|\delta\|^2$ :

$$\|\delta\|^2 = \left\| \frac{1-\alpha_{t-1}}{1-\alpha_t} \cdot \sqrt{1-\beta_t} \cdot z_t + \frac{\beta_t\sqrt{\alpha_{t-1}}}{1-\alpha_t} \cdot x - f_{\phi_t}(z_t) \right\|^2$$

Substitute all terms into the KL divergence expression:

$$D_{KL}[q \parallel p] = \frac{1}{2} \left[ n \log \sigma_t^2 - n \log \left( \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t} \right) - n + \frac{\beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}}{\sigma_t^2} \cdot n + \frac{1}{\sigma_t^2} \cdot \|\delta\|^2 \right]$$

### 1.2.3 (c)

#### Answer

In the KL divergence  $D_{KL}(q(z_{t-1} \mid z_t, x) \parallel P(z_{t-1} \mid z_t; \phi_t))$ :

The term  $q(z_{t-1} \mid z_t, x)$  represents the desired distribution for the denoising step. It is derived from the known data-generating process and incorporates the observed data  $x$  as well as the latent relationship between  $z_t$  and  $z_{t-1}$ . Specifically:

- $q(z_{t-1} \mid z_t, x)$  is a posterior distribution that quantifies the true process of denoising  $z_t$  to obtain  $z_{t-1}$ .
- It is referred to as the target distribution, as it is the distribution we aim to match through the model's predictions.

The term  $P(z_{t-1} \mid z_t; \phi_t)$  represents the approximated distribution predicted by the model during the denoising process. This distribution is parameterized by  $\phi_t$ , which are the learnable parameters of the model. Specifically

- $P(z_{t-1} \mid z_t; \phi_t)$  is the model's attempt to approximate  $q(z_{t-1} \mid z_t, x)$ .
- The parameters  $\phi_t$  are optimized to minimize the divergence between the target  $q(z_{t-1} \mid z_t, x)$  and the prediction  $P(z_{t-1} \mid z_t; \phi_t)$ , thereby improving the accuracy of the denoising steps.