

A. Derivations and proofs in Section 4.3

A.1. Proof of Lemma 1: Approximated expansion of the reconstruction loss

The approximated expansion of the reconstruction loss is mainly the same as Rolínek et al. (2019) except we consider a metric tensor \mathbf{G}_x which is a positive definite Hermitian matrix.

$\delta\check{x}$ and $\delta\hat{x}$ denote $\check{x} - x$ and $\hat{x} - \check{x}$, respectively. Let $\delta z_j \sim \mathcal{N}(\delta z_j; 0, \sigma_{j(x)})$ be an added noise in the reparameterization trick where $z_j = \mu_{j(x)} + \delta z_j$. Then, $\delta\hat{x} = \hat{x} - \check{x}$ is approximated as:

$$\delta\hat{x} \simeq \sum_{j=1}^n \delta z_j \mathbf{x}_{\mu_j}. \quad (26)$$

Next, the reconstruction loss $D(x, \hat{x})$ can be approximated as follows.

$$\begin{aligned} D(x, \hat{x}) &= D(x, x + (\delta\check{x} + \delta\hat{x})) \\ &\simeq {}^t(\delta\check{x} + \delta\hat{x}) \mathbf{G}_x (\delta\check{x} + \delta\hat{x}) \\ &= {}^t\delta\check{x} \mathbf{G}_x \delta\check{x} + {}^t\delta\hat{x} \mathbf{G}_x \delta\hat{x} + 2 {}^t\delta\hat{x} \mathbf{G}_x \delta\check{x} \\ &\simeq D(x, \check{x}) + D(\check{x}, \hat{x}) + \sum_{j=1}^n 2\delta z_j {}^t\mathbf{x}_{\mu_j} \mathbf{G}_x \delta\check{x} \end{aligned} \quad (27)$$

Then, we evaluate the average of $D(x, \hat{x})$ over $z \sim q_\phi(z|x)$, i.e., $\delta z_j \sim \mathcal{N}(\delta z_j; 0, \sigma_{j(x)})$ for all j . Note that $E[\delta z_j \delta z_k] = \sigma_{j(x)}^2 \delta_{jk}$ where δ_{jk} is the Kronecker delta. First, the average of $D(x, \check{x})$ in the last line of Eq. 27 is still $D(x, \check{x})$ since this term does not depend on δz_j . Second, the average of $D(\check{x}, \hat{x})$ in the last line of Eq. 27 is approximated as:

$$\begin{aligned} E_{z \sim q_\phi(z|x)} [D(\check{x}, \hat{x})] &\simeq E_{z \sim q_\phi(z|x)} [{}^t\delta\hat{x} \mathbf{G}_x \delta\hat{x}] \\ &\simeq E_{z \sim q_\phi(z|x)} \left[\left(\sum_{j=1}^n \delta z_j {}^t\mathbf{x}_{\mu_j} \right) \mathbf{G}_x \left(\sum_{k=1}^n \delta z_k \mathbf{x}_{\mu_k} \right) \right] \\ &= \sum_{j=1}^n \sum_{k=1}^n E_{z \sim q_\phi(z|x)} [\delta z_j \delta z_k] {}^t\mathbf{x}_{\mu_j} \mathbf{G}_x \mathbf{x}_{\mu_k} \\ &= \sum_{j=1}^n \sigma_{j(x)}^2 {}^t\mathbf{x}_{\mu_j} \mathbf{G}_x \mathbf{x}_{\mu_j}. \end{aligned} \quad (28)$$

Third, the average of the third term in the last line of Eq. 27, i.e., $\sum_{j=1}^n 2\delta z_j {}^t\mathbf{x}_{\mu_j} \mathbf{G}_x \delta\check{x}$, is 0 since the average of δz_j over $\mathcal{N}(\delta z_j; 0, \sigma_{j(x)})$ is 0.

As a result, the average of $D(x, \hat{x})$ over $z \sim q_\phi(z|x)$ can be approximated as:

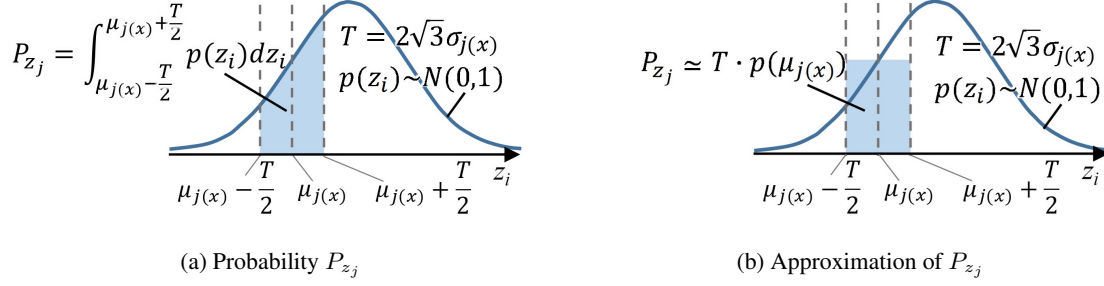
$$E_{z \sim q_\phi(z|x)} [D(x, \hat{x})] \simeq D(x, \check{x}) + \sum_{j=1}^n \sigma_{j(x)}^2 {}^t\mathbf{x}_{\mu_j} \mathbf{G}_x \mathbf{x}_{\mu_j}. \quad (29)$$

A.2. Proof of Lemma 2: More precise derivation of KL divergence approximation.

This appendix explains more precise derivation of KL divergence approximation. First, we show the approximation for the Gaussian prior. We also show at the end of this appendix that our approximation also holds for arbitrary prior.

In the case of Gaussian prior, we show that KL divergence can be interpreted as an amount of information in the transform coding (Goyal, 2001) allowing the distortion $\sigma_{j(x)}^2$. In the transform coding, input data is transformed by an orthonormal transform. Then, the transformed data is quantized, and an entropy code is assigned to the quantized symbol, such that the length of the entropy code is equivalent to the logarithm of the estimated symbol probability. Here, we assume $\sigma_{j(x)}^2 \ll 1$ will be observed in meaningful dimensions as shown later.

It is generally intractable to derive the rate and distortion of individual symbols in the ideal information coding. Thus, we first discuss the case of uniform quantization. Let P_{z_j} and R_{z_j} be the probability and amount of information in the


 Figure 8. Probability for a symbol with mean μ and noise σ^2

uniform quantization coding of $z_j \sim \mathcal{N}(z_j; 0, 1)$. Here, $\mu_{j(x)}$ and $\sigma_{j(x)}^2$ are regarded as a quantized value and a coding noise after the uniform quantization, respectively. Since we assume $\sigma_{j(x)}^2 \ll 1$, $\mu_{j(x)} \sim \mathcal{N}(\mu_{j(x)}; 0, 1)$ will also hold. Let T be a quantization step size. The coding noise after quantization is $T^2/12$ for the quantization step size T , as explained in Appendix G.1. Thus, T is derived as $T = 2\sqrt{3}\sigma_{j(x)}$ from $\sigma_{j(x)}^2 = T^2/12$. We also assume $\sigma_{j(x)}^2 \ll 1$. As shown in Fig.8a, P_{z_j} is denoted by $\int_{\mu_{j(x)}-T/2}^{\mu_{j(x)}+T/2} p(z_j) dz_j$ where $p(z_j)$ is $\mathcal{N}(z_j; 0, 1)$. Using Simpson's numerical integration method and $e^x = 1 + x + O(x^2)$ expansion, P_{z_j} is approximated as:

$$\begin{aligned}
 P_{z_j} &\simeq \frac{T}{6} (p(\mu_{j(x)} - \frac{T}{2}) + 4p(\mu_{j(x)}) + p(\mu_{j(x)} + \frac{T}{2})) \\
 &= \frac{T p(\mu_{j(x)})}{6} \left(4 + e^{\frac{4\mu_{j(x)}T - T^2}{8}} + e^{\frac{-4\mu_{j(x)}T - T^2}{8}} \right) \\
 &\simeq T p(\mu_{j(x)}) (1 - T^2/24) \\
 &= \sqrt{\frac{6}{\pi}} \sigma_{j(x)} e^{-(\mu_{j(x)})^2/2} \left(1 - \frac{\sigma_{j(x)}^2}{2} \right). \tag{30}
 \end{aligned}$$

Using $\log(1+x) = x + O(x^2)$ expansion, $R_{\mu\sigma}$ is derived as:

$$R_{z_j} = -\log P_{z_j} \simeq \frac{1}{2} \left(\mu_{j(x)}^2 + \sigma_{j(x)}^2 - \log \sigma_{j(x)}^2 - \log \frac{6}{\pi} \right) = D_{\text{KL}j(x)}(\cdot) + \frac{1}{2} \log \frac{\pi e}{6}. \tag{31}$$

When R_{z_j} and $D_{\text{KL}j(x)}(\cdot)$ in Eq. 2 are compared, both equations are equivalent except a small constant difference $\frac{1}{2} \log(\pi e/6) \simeq 0.176$ for each dimension. As a result, KL divergence for j -th dimension is equivalent to the rate for the uniform quantization coding, allowing a small constant difference.

To make theoretical analysis easier, we use the simpler approximation as $P_{z_j} = T p(\mu_{j(x)}) = 2\sqrt{3}\sigma_{j(x)} p(\mu_{j(x)})$ instead of Eq.30, as shown in Fig.8b. Then, R_{z_j} is derived as:

$$R_{z_j} = -\log(2\sqrt{3}\sigma_{j(x)} p(\mu_{j(x)})) = \frac{1}{2} (\mu_{j(x)}^2 - \log \sigma_{j(x)}^2 - 1) + \frac{1}{2} \log \frac{\pi e}{6}. \tag{32}$$

Here, the first term of the right equation is equivalent to Eq. 11. This equation also means that the approximation of KL divergence in Eq. 11 is equivalent to the rate in the uniform quantization coding with $P_{z_j} = 2\sqrt{3}\sigma_{j(x)} p(\mu_{j(x)})$ approximation, allowing the same small constant difference as in Eq. 31. It is noted that the approximation $P_{z_j} = 2\sqrt{3}\sigma_{j(x)} p(\mu_{j(x)})$ in Figure 8b can be applied to any kinds of prior PDFs because there is no explicit assumption for the prior PDF. This implies that the theoretical discussion after Eq. 11 in the main text will hold in arbitrary prior PDFs.

The meaning of the small constant difference $\frac{1}{2} \log \frac{\pi e}{6}$ in Eqs. 31 and 32 can be explained as follows: Pearlman & Said (2011) show that the difference of the rate between the ideal information coding and uniform quantization is $\frac{1}{2} \log \frac{\pi e}{6}$. This is caused by the entropy difference of the noise distributions. In the ideal case, the noise distribution is known as a Gaussian. In the case the noise variance is σ^2 , the entropy of the Gaussian noise is $\frac{1}{2} \log(\sigma^2 2\pi e)$. For the uniform quantization with a uniform noise distribution, the entropy is $\frac{1}{2} \log(\sigma^2 12)$. As a result, the difference is just $\frac{1}{2} \log \frac{\pi e}{6}$. Because the rate estimation in this appendix uses a uniform quantization, the small offset $\frac{1}{2} \log \frac{\pi e}{6}$ can be regarded as a difference between

the ideal information coding and the uniform quantization. As a result, KL divergence in Eq. 2 and Eq. 11 can be regarded as a rate in the ideal informaton coding for the symbol with the mean $\mu_{j(\mathbf{x})}$ and variance $\sigma_{j(\mathbf{x})}^2$.

Here, we validate the assumption that $\sigma_{j(\mathbf{x})} \ll 1$ will be observed in meaningful dimensions. From the discussion above, the information R_{z_j} in each dimension can be considered as KL divergence:

$$R_{z_j} = \frac{1}{2} (\mu_{j(\mathbf{x})}^2 + \sigma_{j(\mathbf{x})}^2 - \log \sigma_{j(\mathbf{x})}^2 - 1). \quad (33)$$

For simple analysis, we assume that $\sigma_{j(\mathbf{x})}$ is constant in the j -th dimension. We further assume $\mu_{j(\mathbf{x})} \sim \mathcal{N}(\mu_{j(\mathbf{x})}; 0, 1)$. Then $E[R_{z_j}]$ which shows the information of the j -th dimensional component is derived as:

$$\begin{aligned} E[R_{z_j}] &\simeq E_{\mu_{j(\mathbf{x})} \sim \mathcal{N}(\mu_{j(\mathbf{x})}; 0, 1)}[R_{z_j}] \\ &= \int \frac{1}{2} (\mu_{j(\mathbf{x})}^2 + \sigma_{j(\mathbf{x})}^2 - \log \sigma_{j(\mathbf{x})}^2 - 1) \mathcal{N}(\mu_{j(\mathbf{x})}; 0, 1) d\mu_{j(\mathbf{x})} \\ &= \frac{1}{2} (\sigma_{j(\mathbf{x})}^2 - \log \sigma_{j(\mathbf{x})}^2). \end{aligned} \quad (34)$$

From this equation, we can estimate an amount of information in each dimension from the posterior variance. From this equation, it is derived that if the amount of information $E[R_{z_j}]$ is more than about 1.20 nat or 1.73 bit, $\sigma_{j(\mathbf{x})}^2 < 0.1$ holds. In addition, as the information $E[R_{z_j}]$ is increasing, $\sigma_{j(\mathbf{x})}^2$ becomes exponentially decreasing. As a result, the assumption that $\sigma_{j(\mathbf{x})}^2 \ll 1$ will be observed in meaningful dimensions is reasonable.

Finally, we show that the approximation of the KL divergence in the second line of Eq. 10 also holds for arbitrary priors. Let $p(\mathbf{z})$, $q_\phi(\mathbf{z}|\mathbf{x})$, and $D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$ be an arbitrary prior, a posterior $\prod_j \mathcal{N}(\mu_{j(\mathbf{x})}, \sigma_{j(\mathbf{x})})$, and KL divergence, respectively. First, the shape of $q_\phi(\mathbf{z}|\mathbf{x})$ becomes close to a delta function $\delta(\mathbf{z} - \boldsymbol{\mu}_{(\mathbf{x})})$ when each $\sigma_{j(\mathbf{x})}$ is small. Thus $q_\phi(\mathbf{z}|\mathbf{x})$ will act like a delta function $\delta(\mathbf{z} - \boldsymbol{\mu}_{(\mathbf{x})})$. Next the differential entropy of $q_\phi(\mathbf{z}|\mathbf{x})$, i.e., $-\int q_\phi(\mathbf{z}|\mathbf{x}) \log q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z}$ is derived as $\sum_j^n \log \sigma_{j(\mathbf{x})} \sqrt{2\pi e}$. Using these equations, KL divergence for an arbitrary prior can be approximated by the second line of Eq. 10 as follows:

$$\begin{aligned} D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) &= -\int q_\phi(\mathbf{z}|\mathbf{x}) \log p(\mathbf{z}) d\mathbf{z} + \int q_\phi(\mathbf{z}|\mathbf{x}) \log q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &\simeq -\int \delta(\mathbf{z} - \boldsymbol{\mu}_{(\mathbf{x})}) \log p(\mathbf{z}) d\mathbf{z} + \int q_\phi(\mathbf{z}|\mathbf{x}) \log q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= -\log p(\boldsymbol{\mu}_{(\mathbf{x})}) - \sum_j^n \log \sigma_{j(\mathbf{x})} \sqrt{2\pi e} \\ &= -\log \left(p(\boldsymbol{\mu}_{(\mathbf{x})}) \prod_{j=1}^n \sigma_{j(\mathbf{x})} \right) - \frac{n \log 2\pi e}{2} \end{aligned} \quad (35)$$

In the derivation of the third line in Eq. 10, we assume that $q_\phi(\boldsymbol{\mu}_{(\mathbf{x})})$ is close to $p(\boldsymbol{\mu}_{(\mathbf{x})})$ where $p(\cdot)$ is the prior distribution of \mathbf{z} . The reason of this assumption is as follows: ELBO can be also derived as $\log p(\mathbf{x}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$ (Bishop, 2006). When ELBO is maximized at each \mathbf{x} , $q_\phi(\mathbf{z}|\mathbf{x}) \simeq p_\theta(\mathbf{z}|\mathbf{x})$ will hold to minimize KL divergence where $\log p(\mathbf{x})$ is a constant. Finally, we have $q_\phi(\mathbf{z}) \simeq p(\mathbf{z})$ by the marginalization of \mathbf{x} . Appendix A.3 also validates this assumption in the simple 1-dimensional VAE case where $q_\phi(\mu_{(x)}) = p(\mu_{(x)}) = \mathcal{N}(\mu_{(x)}; 0, 1)$ holds. Thus, we can derive the approximation $p(\boldsymbol{\mu}_{(\mathbf{x})}) \simeq q_\phi(\boldsymbol{\mu}_{(\mathbf{x})}) = p(\mathbf{x}) |\det(\partial \mathbf{x} / \partial \boldsymbol{\mu}_{(\mathbf{x})})|$ in Eq. 10.

A.3. Proof of Lemma 3: Estimation of the coding loss and transform loss in 1-dimensional linear VAE

This appendix estimates the coding loss and transform loss in 1-dimensional linear β -VAE for the Gaussian data, and also shows that the result is consistent with the Wiener filter (Wiener, 1964). Let x be a one dimensional input data with the normal distribution:

$$x \in \mathbb{R}, \quad x \sim \mathcal{N}(x; 0, \sigma_x^2). \quad (36)$$

First, a simple VAE model in this analysis is explained. z denotes a one dimensional latent variable. Let the prior distribution $p(z)$ be $\mathcal{N}(z; 0, 1)$. Next, two linear parametric encoder and decoder are provided with constant parameters a , b , and σ_z to

optimize:

$$\begin{aligned} \text{Enc}_\phi : \quad z &= \mu + \sigma_z \epsilon \text{ where } \mu = ax \text{ and } \epsilon \sim \mathcal{N}(\epsilon; 0, 1), \\ \text{Dec}_\theta : \quad \hat{x} &= bz. \end{aligned} \quad (37)$$

Here, the encoding parameter ϕ consists of $\{a, \sigma_z\}$, and the decoding parameter θ consists of $\{b\}$. Then the square error is used as a reconstruction loss.

Next, the objective is derived. $D_{\text{KL}x}$ and D_x denote the KL divergence and reconstruction loss at x , respectively. We further assume that D_x uses a square error. Then we define the loss objective at x as $L_x = D_x + \beta D_{\text{KL}x}$. Using Eq. 11, $D_{\text{KL}x}$ can be evaluated as:

$$\begin{aligned} D_{\text{KL}x} &= -\log(\sigma_z p(\mu)) - \frac{1}{2} \log 2\pi e \\ &= -\log(\sigma_z \mathcal{N}(ax; 0, 1)) - \frac{1}{2} \log 2\pi e \\ &= -\log \sigma_z + \frac{a^2 x^2}{2} - \frac{1}{2}. \end{aligned} \quad (38)$$

Then, D_x is evaluated as:

$$\begin{aligned} D_x &= E_{\epsilon \sim \mathcal{N}(\epsilon; 0, 1)} \left[(x - \text{Dec}_\theta(\text{Enc}_\phi(x)))^2 \right] \\ &= E_{\epsilon \sim \mathcal{N}(\epsilon; 0, 1)} \left[(x - (b(ax + \sigma_z \epsilon)))^2 \right] \\ &= \int ((ab - 1)^2 x^2 + 2(ab - 1)xb\sigma_z \epsilon + b^2 \sigma_z^2 \epsilon^2) \mathcal{N}(\epsilon; 0, 1) d\epsilon \\ &= (ab - 1)^2 x^2 + b^2 \sigma_z^2. \end{aligned} \quad (39)$$

By averaging L_x over $x \sim \mathcal{N}(x; 0, \sigma_x^2)$, the objective L to minimize is derived as:

$$\begin{aligned} L &= E_{x \sim \mathcal{N}(x; 0, \sigma_x^2)} [L_x] \\ &= \int \left((ab - 1)^2 x^2 + b^2 \sigma_z^2 + \beta \left(-\log \sigma_z + \frac{a^2 x^2}{2} - \frac{1}{2} \right) \right) \mathcal{N}(x; 0, \sigma_x^2) dx. \\ &= (ab - 1)^2 \sigma_x^2 + b^2 \sigma_z^2 + \beta \left(-\log \sigma_z + \frac{a^2 \sigma_x^2}{2} - \frac{1}{2} \right). \end{aligned} \quad (40)$$

Here, the first term $(ab - 1)^2 \sigma_x^2$ and the second term $b^2 \sigma_z^2$ in the last line are corresponding to the transform loss D_T and coding loss D_C , respectively.

By solving $dL/da = 0$, $dL/db = 0$, and $dL/d\sigma_z = 0$, a , b , and σ_z are derived as follows:

$$\begin{aligned} a &= 1/\sigma_x, \\ b &= \frac{\sigma_x \left(1 + \sqrt{1 - 2\beta/\sigma_x^2} \right)}{2}, \\ \sigma_z &= \frac{2\sqrt{\beta/2}}{\sigma_x \left(1 + \sqrt{1 - 2\beta/\sigma_x^2} \right)}. \end{aligned} \quad (41)$$

From Eq. 41, D_T and D_C are derived as:

$$\begin{aligned} D_T &= \left(\frac{\sqrt{1 - 2\beta/\sigma_x^2} - 1}{2} \right)^2 \sigma_x^2, \\ D_C &= \beta/2. \end{aligned} \quad (42)$$

As shown in section 4.1, the added noise, $\beta/2$, should be reasonably smaller than the data variance σ_x^2 . If $\sigma_x^2 \gg \beta$, b and σ_z in Eq. 41 can be approximated as:

$$D_T \simeq \frac{(\beta/2)^2}{\sigma_x^2} = \frac{\beta/2}{\sigma_x^2} D_C. \quad (43)$$

As shown in this equation, D_T/D_C is small in the VAE where the added noise is reasonably small, and D_T can be ignored.

Note that the distribution of $\mu = a x = x/\sigma_x$, i.e., $q_\phi(\mu)$, is derived as $\mathcal{N}(\mu; 0, 1)$ by scaling $p(x) = \mathcal{N}(x; 0, \sigma_x^2)$ with a factor of $a = 1/\sigma_x$. Thus $q_\phi(\mu)$ is equivalent to the prior of z , i.e., $\mathcal{N}(z; 0, 1)$ in this simple VAE case.

Next, the relation to the Wiener filter (Wiener, 1964) is discussed. The Wiener filter is one of the most basic, but most important theories for signal restoration. We consider an simple 1-dimensional Gaussian process. Let $x \sim \mathcal{N}(x; 0, \sigma_x^2)$ be input data. Then, x is scaled by s , and a Gaussian noise $n \sim \mathcal{N}(n; 0, \sigma_n^2)$ is added. Thus, $y = s x + n$ is observed. From the Wiener filter theory, the estimated value with minimum distortion, \hat{x} can be formulated as:

$$\hat{x} = \frac{s\sigma_x^2}{s^2\sigma_x^2 + \sigma_n^2} y. \quad (44)$$

In this case, the estimation error is derived as:

$$E[(\hat{x} - x)^2] = \frac{\sigma_n^4}{(s^2\sigma_x^2 + \sigma_n^2)^2} \sigma_x^2 + \frac{s^2\sigma_x^4}{(s^2\sigma_x^2 + \sigma_n^2)^2} \sigma_n^2 = \frac{\sigma_x^2}{\sigma_x^2 + (\sigma_n^2/s^2)} (\sigma_n^2/s^2). \quad (45)$$

In the second equation, the first term is corresponding to the transform loss, and the second term is corresponding to the coding loss. Here the ratio of the transform loss and coding loss is derived as $\sigma_n^2/(s^2\sigma_x^2)$. By appying $s = 1/\sigma_x$ and $\sigma_n = \sigma_z$ to $\sigma_n^2/(s^2\sigma_x^2)$ and assuming $\sigma_x^2 \gg \beta/2$, this ratio can be described as:

$$\frac{\sigma_n^2}{s^2\sigma_x^2} = \sigma_z^2 = \frac{\beta/2}{\sigma_x^2} \frac{4}{\left(1 + \sqrt{1 - 2\beta/\sigma_x^2}\right)^2} = \frac{\beta/2}{\sigma_x^2} + O\left(\left(\frac{\beta/2}{\sigma_x^2}\right)^2\right). \quad (46)$$

This result is consistent with Eq. 43, implying that optimized VAE and the Wiener filter show similar behaviours.

A.4. Proof of Lemma 4 : Derivation of the orthogonality

Lemma 4 is proved by examining the minimum condition of L_x at x . The proof outline is similar to Kato et al. (2020) while $\sigma_{j(x)}$ should be also considered as a variable in our derivation.

We first show the following mathematical formula which is used our derivation. Let \mathbf{A} be a regular matrix and \mathbf{a}_i be its i -th column vector. $\tilde{\mathbf{a}}_i$ denotes the i -th column vector of a cofactor matrix for \mathbf{A} . Then the following equation holds mathematically.

$$\frac{d \log |\det(\mathbf{A})|}{d\mathbf{a}_i} = \frac{d \log |\det(\mathbf{A})|}{d \det(\mathbf{A})} \frac{d \det(\mathbf{A})}{d\mathbf{a}_i} = \frac{1}{\det(\mathbf{A})} \tilde{\mathbf{a}}_i. \quad (47)$$

Let $\tilde{\mathbf{x}}_{\mu_j}$ be the j -th column vector of a cofactor matrix for Jacobian matrix $\partial \mathbf{x} / \partial \boldsymbol{\mu}_{(x)}$. Using the formula in Eq. 47, the partial derivative of L_x by \mathbf{x}_{μ_j} is described by

$$\frac{\partial L_x}{\partial \mathbf{x}_{\mu_j}} = 2\sigma_{j(x)}^2 \mathbf{G}_x \mathbf{x}_{\mu_j} - \frac{\beta}{\det(\partial \mathbf{x} / \partial \boldsymbol{\mu}_{(x)})} \tilde{\mathbf{x}}_{\mu_j}. \quad (48)$$

Note that ${}^t \mathbf{x}_{\mu_k} \cdot \tilde{\mathbf{x}}_{\mu_j} = \det(\partial \mathbf{x} / \partial \mathbf{z}) \delta_{jk}$ holds by the cofactor's property. Here, \cdot denotes the dot product, and δ_{jk} denotes the Kronecker delta. By setting Eq. 48 to zero and multiplying ${}^t \mathbf{x}_{\mu_k}$ from the left, we have the next orthogonal form of \mathbf{x}_{μ_j} :

$$(2\sigma_{j(x)}^2 / \beta) {}^t \mathbf{x}_{\mu_k} \mathbf{G}_x \mathbf{x}_{\mu_j} = \delta_{jk}. \quad (49)$$

Next, the partial derivative of L_x by $\sigma_{j(x)}$ is derived as:

$$\frac{\partial L_x}{\partial \sigma_{j(x)}} = 2\sigma_{j(x)} {}^t \mathbf{x}_{\mu_j} \mathbf{G}_x \mathbf{x}_{\mu_j} - \frac{\beta}{\sigma_{j(x)}}. \quad (50)$$

By setting Eq. 50 to zero, we have the next equation:

$$(2\sigma_{j(x)}^2 / \beta) {}^t \mathbf{x}_{\mu_j} \mathbf{G}_x \mathbf{x}_{\mu_j} = 1. \quad (51)$$

Note that Eq. 51 is a part of Eq. 49 where $j = k$. As a result, the condition to minimize L_x is derived as Eq. 49.

A.5. Proof of Proposition 1: Estimation of input data distribution in the metric space

This equation explains the derivation of Eq. 19 in Proposition 1. Using Eq. 16, the third equation in Eq. 19 is derived as:

$$p(\mathbf{y}) = \prod_j^n p(y_j) = \prod_j^n (dy_j/d\mu_{j(\mathbf{x})})^{-1} p(\mu_j) = \prod_j^n p(\mu_j) \prod_j^n \frac{\sigma_{j(\mathbf{x})}}{\sqrt{\beta/2}} = (\beta/2)^{n/2} p(\mu_{j(\mathbf{x})}) \prod_j^n \sigma_{j(\mathbf{x})}. \quad (52)$$

This shows that the posterior variance $\sigma_{j(\mathbf{x})}$ bridges between the distributions of data and prior. Thus the prior close to the data distribution will facilitate training, where $\sigma_{j(\mathbf{x})}$ is close to constant.

The fourth equation in Eq. 19 in Proposition 1 is derived by applying Eq. 16 to Eq. 13 and arranging the result. Let $L_{\min \mathbf{x}}$ be a minimum of $L_{\mathbf{x}}$ at \mathbf{x} . $D_{\min \mathbf{x}}$ and $R_{\min \mathbf{x}}$ denote a coding loss and KL divergence in $L_{\min \mathbf{x}}$, respectively.

First, $D_{\min \mathbf{x}}$ is derived. The next equation holds from Eq. 14.

$$\sigma_{j(\mathbf{x})}^2 {}^t \mathbf{x}_{\mu_j} \mathbf{G}_{\mathbf{x}} \mathbf{x}_{\mu_j} = \beta/2. \quad (53)$$

By applying Eq. 53 to the first term of Eq. 13, $D_{\min \mathbf{x}}$ is derived as:

$$D_{\min \mathbf{x}} = \sum_j^n \sigma_{j(\mathbf{x})}^2 {}^t \mathbf{x}_{\mu_j} \mathbf{G}_{\mathbf{x}} \mathbf{x}_{\mu_j} = n\beta/2. \quad (54)$$

This implies that the reconstruction loss is constant for all inputs at the minimum condition.

Second, $R_{\min \mathbf{x}}$ is derived. From Eq. 52, the next equation holds.

$$p(\mu_{j(\mathbf{x})}) \prod_j^n \sigma_{j(\mathbf{x})} = (\beta/2)^{-n/2} p(\mathbf{y}). \quad (55)$$

By applying Eq. 55 to the second equation of Eq. 10, $R_{\min \mathbf{x}}$ is derived as:

$$R_{\min \mathbf{x}} = -\log p(\mathbf{y}) - \frac{n \log(\beta\pi e)}{2}. \quad (56)$$

As a result, the minimum value of the objective $L_{\min \mathbf{x}}$ is derived as:

$$L_{\min \mathbf{x}} = D_{\min \mathbf{x}} + \beta R_{\min \mathbf{x}} = -\beta \log p(\mathbf{y}) + \frac{n\beta}{2} (1 - \log(\beta\pi e)). \quad (57)$$

As a result, $p(\mathbf{x})$ can be evaluated as:

$$\exp(-L_{\min \mathbf{x}}/\beta) = p(\mathbf{y}) \exp(-\frac{n(1 - \log(\beta\pi e))}{2}) \propto p(\mathbf{y}) \simeq p_{\mathbf{G}_{\mathbf{x}}}(\mathbf{x}). \quad (58)$$

This result implies that the VAE objective converges to the log-likelihood of the input \mathbf{x} at the optimized condition as expected.

A.6. Proof of Proposition 2: Estimation of data distribution in the input space

This appendix shows the derivation of variables in Eqs. 19 and 21. When we estimate a probability in real dataset, we use an approximation of $L_{\mathbf{x}}$. First, the derivation of $L_{\mathbf{x}}$ approximation for the input \mathbf{x} is presented. Then, the PDF ratio between the input space and inner product space is explained for the cases $m = n$ and $m > n$.

Derivation of $L_{\mathbf{x}}$ approximation for the input \mathbf{x} of real data:

As shown in in Eq. 1, $L_{\mathbf{x}}$ is denoted as $-E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})}[\cdot] + \beta D_{\text{KL}}(\cdot)$. We approximate $E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})}[\cdot]$ as $\frac{1}{2}(D(\mathbf{x}, \text{Dec}_{\theta}(\mu_{\mathbf{x}} + \sigma_{\mathbf{x}})) + D(\mathbf{x}, \text{Dec}_{\theta}(\mu_{\mathbf{x}} - \sigma_{\mathbf{x}})))$, i.e., the average of two samples, instead of the average over $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. $D_{\text{KL}}(\cdot)$ can be calculated from $\mu_{\mathbf{x}}$ and $\sigma_{\mathbf{x}}$ using Eq. 2.

The PDF ratio in the case $m = n$:

The PDF ratio for $m = n$ is a Jacobian determinant between two spaces. First, $(\frac{\partial \mathbf{x}}{\partial \mathbf{y}})^T \mathbf{G}_{\mathbf{x}} (\frac{\partial \mathbf{x}}{\partial \mathbf{y}}) = \mathbf{I}_m$ holds from

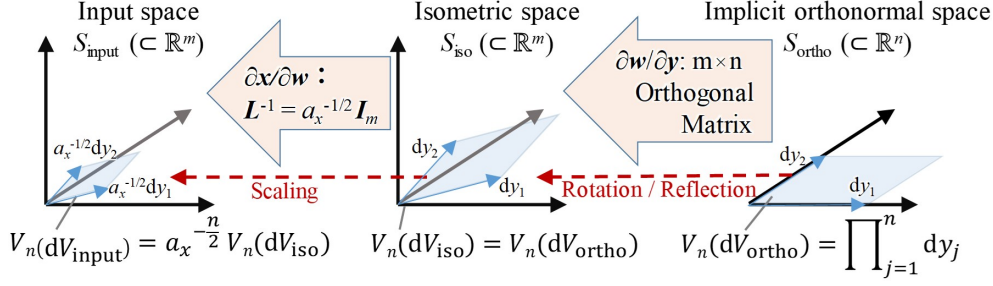


Figure 9. Projection of the volume element from the implicit orthonormal space to the isometric space and input space. $V_n(\cdot)$ denotes n -dimensional volume.

Eq. 17. $|\partial \mathbf{x} / \partial \mathbf{y}|^2 |\mathbf{G}_x| = 1$ also holds by calculating the determinant. Finally, $|\partial \mathbf{x} / \partial \mathbf{y}|$ is derived as $|\mathbf{G}_x|^{1/2}$ using $|\partial \mathbf{y} / \partial \mathbf{x}| = |\partial \mathbf{x} / \partial \mathbf{y}|^{-1}$.

The PDF ratio in the case $m > n$ and $\mathbf{G}_x = a_x \mathbf{I}_m$:

Although the strict derivation needs the treatment of the Riemannian manifold, we provide a simple explanation in this appendix. Here, it is assumed that $D_{\text{KL}(j)}(\cdot) > 0$ holds for all $j = [1, ..n]$. If $D_{\text{KL}(j)}(\cdot) = 0$ for some j , n is replaced by the number of latent variables with $D_{\text{KL}(j)}(\cdot) > 0$.

For the implicit isometric space $S_{\text{iso}} (\subset \mathbb{R}^m)$, there exists a matrix \mathbf{L}_x such that both $\mathbf{y} = \mathbf{L}_x \mathbf{x}$ and $\mathbf{G}_x = {}^t \mathbf{L}_x \mathbf{L}_x$ holds. \mathbf{w} denotes a point in S_{iso} , i.e., $\mathbf{w} \in S_{\text{iso}}$. Because \mathbf{G}_x is assumed as $a_x \mathbf{I}_m$ in Section 4.3, $\mathbf{L}_x = a_x^{1/2} \mathbf{I}_m$ holds. Then, the mapping function $\mathbf{w} = h(\mathbf{x})$ between S_{input} and S_{iso} is defined, such that:

$$\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \mathbf{w}}{\partial \mathbf{x}} = \mathbf{L}_x, \text{ and } h(\mathbf{x}^{(0)}) = \mathbf{w}^{(0)} \text{ for } \exists \mathbf{x}^{(0)} \in S_{\text{input}} \text{ and } \exists \mathbf{w}^{(0)} \in S_{\text{iso}}. \quad (59)$$

Let $\delta \mathbf{x}$ and $\delta \mathbf{w}$ are infinitesimal displacements around \mathbf{x} and $\mathbf{w} = h(\mathbf{x})$, such that $\mathbf{w} + \delta \mathbf{w} = h(\mathbf{x} + \delta \mathbf{x})$. Then the next equation holds from Eq. 59:

$$\delta \mathbf{w} = \mathbf{L}_x \delta \mathbf{x}. \quad (60)$$

Let $\delta \mathbf{x}^{(1)}$, $\delta \mathbf{x}^{(2)}$, $\delta \mathbf{w}^{(1)}$, and $\delta \mathbf{w}^{(2)}$ be two arbitrary infinitesimal displacements around \mathbf{x} and $\mathbf{w} = h(\mathbf{x})$, such that $\delta \mathbf{w}^{(1)} = \mathbf{L}_x \delta \mathbf{x}^{(1)}$ and $\delta \mathbf{w}^{(2)} = \mathbf{L}_x \delta \mathbf{x}^{(2)}$. Then the following equation holds, where \cdot denotes the dot product.

$${}^t \delta \mathbf{x}^{(1)} \mathbf{G}_x \delta \mathbf{x}^{(2)} = {}^t (\mathbf{L}_x \delta \mathbf{x}^{(1)}) (\mathbf{L}_x \delta \mathbf{x}^{(2)}) = \delta \mathbf{w}^{(1)} \cdot \delta \mathbf{w}^{(2)}. \quad (61)$$

This equation shows the isometric mapping from the inner product space for $\mathbf{x} \in S_{\text{input}}$ with the metric tensor \mathbf{G}_x to the Euclidean space for $\mathbf{w} \in S_{\text{iso}}$.

Note that all of the column vectors in the Jacobian matrix $\partial \mathbf{x} / \partial \mathbf{y}$ also have a unit norm and are orthogonal to each other in the metric space for $\mathbf{x} \in S_{\text{input}}$ with the metric tensor \mathbf{G}_x . Therefore, the $m \times n$ Jacobian matrix $\partial \mathbf{w} / \partial \mathbf{y}$ should have a property that all of the column vectors have a unit norm and are orthogonal to each other in the Euclidean space.

Then n -dimensional space which is composed of the meaningful dimensions from the implicit isometric space is named as the implicit orthonormal space S_{ortho} . Figure 9 shows the projection of the volume element from the implicit orthonormal space to the isometric space and input space. Let dV_{ortho} be an infinitesimal n -dimensional volume element in S_{ortho} . This volume element is a n -dimensional rectangular solid having each edge length dy_j . Let $V_n(dV_{\text{ortho}})$ be the n -dimensional volume of a volume element dV_{ortho} . Then, $V_n(dV_{\text{ortho}}) = \prod_j^n dy_j$ holds. Next, dV_{ortho} is projected to n dimensional infinitesimal element dV_{iso} in S_{iso} by $\partial \mathbf{w} / \partial \mathbf{y}$. Because of the orthonormality, dV_{iso} is equivalent to the rotation / reflection of dV_{ortho} , and $V_n(dV_{\text{iso}})$ is the same as $V_n(dV_{\text{ortho}})$, i.e., $\prod_j^n dy_j$. Then, dV_{iso} is projected to n -dimensional element dV_{input} in S_{input} by $\partial \mathbf{x} / \partial \mathbf{w} = \mathbf{L}_x^{-1} = a_x^{-1/2} \mathbf{I}_m$. Because each dimension is scaled equally by the scale factor $a_x^{-1/2}$, $V_n(dV_{\text{input}}) = \prod_j^n a_x^{-1/2} dy_j = a_x^{-n/2} V_n(dV_{\text{ortho}})$ holds. Here, the ratio of the volume element between S_{input} and S_{ortho} is $V_n(dV_{\text{input}}) / V_n(dV_{\text{ortho}}) = a_x^{-n/2}$. Note that the PDF ratio is derived by the reciprocal of $V_n(dV_{\text{input}}) / V_n(dV_{\text{ortho}})$. As a result, the PDF ratio is derived as $a_x^{n/2}$.

A.7. Proof of proposition 3: Determination of the meaningful dimension for representation

This appendix explain the derivation of Proposition 3. Here, we estimate the KL divergence, i.e., a rate for the dimensions whose variance is less than β . As shown later, the discussion in this appendix is closely related with Rate-distortion theory (Berger, 1971; Pearlman & Said, 2011; Goyal, 2001).

Let L_G , D_G , and R_G be averages of $L_{\min \mathbf{x}}$, $D_{\min \mathbf{x}}$, and $R_{\min \mathbf{x}}$ in Appendix A.5 over $\mathbf{x} \sim p(\mathbf{x})$, respectively. Here, $L_G = D_G + \beta R_G$ holds by definition. Since $D_{\min \mathbf{x}}$ is a constant $n\beta/2$ as in Eq. 54, D_G is derived as:

$$D_G = n\beta/2. \quad (62)$$

As D_G is constant, the minimum condition of L_G is equivalent to that of R_G . Let $D_{\text{KLmin } j(x)}$ be a KL divergence of the j -th dimensional component at the minimum condition. Here, $R_{\min \mathbf{x}} = \sum_j^n D_{\text{KLmin } j(x)}$ holds by definition. Eq. 11 holds for for small $\beta/2$. Thus, we can approximate $D_{\text{KLmin } j(x)}$ for small $\beta/2$ from Eqs. 16 and 11 as:

$$\begin{aligned} D_{\text{KLmin } j(x)} &\simeq -\log(\sigma_{j(\mathbf{x})} p(\mu_{j(\mathbf{x})})) - \frac{\log 2\pi e}{2} \\ &= -\log(\sqrt{\beta/2} p(y_j)) - \frac{\log 2\pi e}{2} \\ &= -\log(p(y_j)) - \frac{\log \beta \pi e}{2} \\ &= -\log(p(y_j)) - H(\mathcal{N}(y_j; 0, \beta/2)). \end{aligned} \quad (63)$$

Here, $H(\mathcal{N}(y_j; 0, \beta/2))$ denotes a entropy of the Gaussian with variance $\beta/2$. Next, R_G is expressed as:

$$\begin{aligned} R_G &= E_{\mathbf{x} \sim p(\mathbf{x})}[R_{\min \mathbf{x}}] \\ &= E_{\mathbf{x} \sim p(\mathbf{x})} \left[\sum_j^n D_{\text{KLmin } j(x)} \right] \\ &\simeq - \int p(\mathbf{x}) \sum_j^n (-\log(p(y_j)) - H(\mathcal{N}(y_j; 0, \beta/2)), 0) d\mathbf{x}. \\ &= - \int p(\mathbf{y}) \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right|^{-1} \sum_j^n (-\log(p(y_j)) - H(\mathcal{N}(y_j; 0, \beta/2)), 0) \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right| d\mathbf{y} \\ &= - \int p(\mathbf{y}) \sum_j^n (-\log(p(y_j)) - H(\mathcal{N}(y_j; 0, \beta/2)), 0) d\mathbf{y}. \\ &= \sum_j^n \left(- \int p(y_j) \log(p(y_j)) dy_j - H(\mathcal{N}(y_j; 0, \beta/2)) \right). \end{aligned} \quad (64)$$

Note that the KL-divergence is always equal or greater than 0 by definition. By considering this, R_G is further approximated as:

$$R_G \simeq \sum_j^n \max \left(- \int p(y_j) \log(p(y_j)) dy_j - H(\mathcal{N}(y_j; 0, \beta/2)), 0 \right). \quad (65)$$

Note that the approximation of Eq. 65 is reasonable from the Rate-distortion theory and optimal transform coding theory (Berger, 1971; Pearlman & Said, 2011; Goyal, 2001). The outline of Rate-distortion theory and optimal transform coding is explained in Appendix B.6. The term $-\int p(y_j) \log p(y_j) dy_j$ is the entropy of y_j . Thus, the optimal implicit isometric space is derived such that the entropy of data representation is minimum. When the data manifold has a disentangled property in the given metric by nature, each y_j will capture a disentangled feature with minimum entropy such that the mutual information between implicit isometric components becomes minimized. This is analogous to PCA for Gaussian data, which gives the disentangled representation with minimum entropy in SSE. Considering the similarity to the PCA eigenvalues, the variance of y_j will indicate the importance of each dimension.

Thus, if the entropy of y_j is larger than $H(\mathcal{N}(0, \beta/2))$, then it is reasonable that $D_{\text{KLmin } j(x)} > 0$ holds. By contrast, if the entropy of y_j is less than $H(\mathcal{N}(0, \beta/2))$, then $D_{\text{KLmin } j(x)} = 0$ will hold. In such dimensions, $\sigma_{j(x)} = 1$, $\mu_{j(x)} = 0$, and $D_{\text{KL}j(x)} = 0$ will hold. In addition, $\sigma_{j(x)}^2 \mathbf{x}_{\mu_{j(x)}}^T \mathbf{G}_{\mathbf{x}} \mathbf{x}_{\mu_{j(x)}}$ will be close to 0 because this needs not to be balanced with $D_{\text{KL}j(x)}$.

These properties of VAE can be clearly explained by rate-distortion theory (Berger, 1971), which has been successfully applied to transform coding such as image / audio compression. Appendix B.6 explains that VAE can be interpreted as an optimal transform coding with non-linear scaling of latent space.

Thus, latent variables with variances from the largest to the n -th with $D_{\text{KL}j(x)} > 0$ are sufficient for the representation and the dimensions with $D_{\text{KL}j(x)} = 0$ can be ignored, allowing the reduction of the dimension n for \mathbf{z} .

A.8. Proof of proposition 4: Derivation of the estimated variance

This appendix explains the derivation of quantitative importance for each dimension in Eq. 22 of Proposition 4.

First, we set y_j to 0 at $\mu_{j(x)} = 0$ to derive y_j value from $dy_j/d\mu_{j(x)}$ in Eq. 16. We also assume that the prior distribution is $\mathcal{N}(\mathbf{z}; 0, \mathbf{I}_n)$. The variance is derived by the subtraction of $E[y_j]^2$, i.e., the square of the mean, from $E[y_j^2]$, i.e., the square mean. Thus, the approximations of both $E[y_j]$ and $E[y_j^2]$ are needed.

First, the approximation of the mean $E[y_j]$ is explained. Because the cumulative distribution functions (CDFs) of y_j are the same as CDF of $\mu_{j(x)}$, the following equations hold:

$$\int_{-\infty}^0 p(y_j) dy_j = \int_{-\infty}^0 p(\mu_{j(x)}) d\mu_{j(x)} = 0.5, \quad \int_0^{\infty} p(y_j) dy_j = \int_0^{\infty} p(\mu_{j(x)}) d\mu_{j(x)} = 0.5. \quad (66)$$

This equation means that the median of the y_j distribution is 0. Because the mean and median are close in most cases, the mean $E[y_j]$ can be approximated as 0. As a result, the variance of y_j can be approximated by the square mean $E[y_j^2]$.

Second, the approximation of the square mean $E[y_j^2]$ is explained. Since we assume the manifold has a disentangled property by nature, the standard deviation of the posterior $\sigma_{j(x)}$ is assumed as a function of $\mu_{j(x)}$, regardless of \mathbf{x} . This function is denoted as $\sigma_j(\mu_{j(x)})$. For ≥ 0 , y_j is approximated as follows, using Eq. 16 and replacing the average of $1/\sigma_j(\dot{\mu}_{j(x)})$ over $\dot{\mu}_{j(x)} = [0, \mu_{j(x)}]$ by $1/\sigma_j(\mu_{j(x)})$:

$$y_j = \int_0^{\mu_{j(x)}} \frac{dy_j}{dz_j} dz_j = \sqrt{\frac{\beta}{2}} \int_0^{z_i} \frac{1}{\sigma_j(\dot{\mu}_{j(x)})} dz_i \simeq \sqrt{\frac{\beta}{2}} \frac{1}{\sigma_j(\mu_{j(x)})} \int_0^{\mu_{j(x)}} dz_j = \sqrt{\frac{\beta}{2}} \frac{\mu_{j(x)}}{\sigma_j(\mu_{j(x)})}. \quad (67)$$

The same approximation is applied to $z_i < 0$. Then the square mean of y_i is approximated as follows, assuming that the correlation between $\sigma(\mu_{j(x)})^{-2}$ and $\mu_{j(x)}^2$ is low:

$$\int y_j^2 p(y_j) dy_j \simeq \frac{\beta}{2} \int \left(\frac{\mu_{j(x)}}{\sigma_j(\mu_{j(x)})} \right)^2 p(\mu_{j(x)}) d\mu_{j(x)} \simeq \frac{\beta}{2} \int \sigma_j(\mu_{j(x)})^{-2} p(\mu_{j(x)}) d\mu_{j(x)} \int \mu_{j(x)}^2 p(\mu_{j(x)}) d\mu_{j(x)}. \quad (68)$$

Finally, the square mean of y_i is approximated as the following equation, using $\int \mu_{j(x)}^2 p(\mu_{j(x)}) d\mu_{j(x)} = 1$ and replacing $\sigma_j(\mu_{j(x)})^2$ by $\sigma_{j(x)}^2$, i.e., the posterior variance derived from the input data:

$$\int y_j^2 p(y_j) dy_j \simeq \frac{\beta}{2} \int \sigma_j(\mu_{j(x)})^{-2} p(\mu_{j(x)}) d\mu_{j(x)} \simeq \frac{\beta}{2} \int \frac{E}{\mu_{j(x)} \sim p(\mu_{j(x)})} [\sigma_j(\mu_{j(x)})^{-2}] \simeq \frac{\beta}{2} \int \frac{E}{\mathbf{x} \sim p(\mathbf{x})} [\sigma_{j(x)}^{-2}]. \quad (69)$$

Although some rough approximations are used in the expansion, the estimated variance in the last equation seems still reasonable, because $\sigma_{j(x)}$ shows a scale factor between y_j and $\mu_{j(x)}$ while the variance of $\mu_{j(x)}$ is always 1 for the prior $\mathcal{N}(\mu_{j(x)}; 0, 1)$. Considering the variance of the prior $\int \mu_{j(x)}^2 p(\mu_{j(x)}) d\mu_{j(x)}$ in the expansion, this estimation method can be applied to any prior distribution.

B. Detailed relation to prior works

This section first describes the clear formulation of ELBO in VAE by utilizing isometric embedding. Then the detailed relationship, including correction, with previous works are explained.

B.1. Derivation of ELBO with clear and quantitative form

This section clarifies that the ELBO value after optimization becomes close to the log-likelihood of input data in the metric space (not input space), by the theoretical derivation of the reconstruction loss and KL divergence via isometric embedding.

We derive the ELBO (without β) at \mathbf{x} in Eq. 1, i.e., $E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$ when the objective of β -VAE $L_{\mathbf{x}}$ (with β) in Eq. 57, i.e., $L_{\mathbf{x}} = D(\mathbf{x}, \hat{\mathbf{x}}) + \beta D_{\text{KL}}(\cdot)$ is optimised.

First, the reconstruction loss can be rewritten as:

$$E_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] = \int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z} = \int q_\phi(\mathbf{y}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{y}) d\mathbf{y}. \quad (70)$$

Let $\mu_{\mathbf{y}(\mathbf{x})}$ be a implicit isometric variable corresponding to $\mu_{(\mathbf{x})}$. Because the posterior variance in each isometric latent variable is a constant $\beta/2$, $q_\phi(\mathbf{y}|\mathbf{x}) \simeq \mathcal{N}(\mathbf{y}; \mu_{\mathbf{y}(\mathbf{x})}, (\beta/2)\mathbf{I}_n)$ will hold. If $\beta/2$ is small, $p(\hat{\mathbf{x}}) \simeq p(\mathbf{x})$ will hold. Then, the next equation will hold also using isometricity;

$$p_\theta(\mathbf{x}|\mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{y}) = p_\theta(\mathbf{x}|\hat{\mathbf{x}}) = p(\hat{\mathbf{x}}|\mathbf{x})p(\mathbf{x})/p(\hat{\mathbf{x}}) \simeq p(\hat{\mathbf{x}}|\mathbf{x}) \simeq q_\phi(\mathbf{y}|\mathbf{x}) \simeq \mathcal{N}(\mathbf{y}; \mu_{\mathbf{y}(\mathbf{x})}, (\beta/2)\mathbf{I}_n). \quad (71)$$

Thus the reconstruction loss is estimated as:

$$\begin{aligned} E_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] &\simeq \int \mathcal{N}(\mathbf{y}; \mu_{\mathbf{y}(\mathbf{x})}, (\beta/2)\mathbf{I}_n) \log \mathcal{N}(\mathbf{y}; \mu_{\mathbf{y}(\mathbf{x})}, (\beta/2)\mathbf{I}_n) d\mathbf{y} \\ &= -(n/2) \log(\beta\pi e). \end{aligned} \quad (72)$$

Next, KL divergence is derived from Eq. 56 as:

$$D_{\text{KL}}(\cdot) = R_{\min \mathbf{x}} = -\log p(\mathbf{y}) - (n/2) \log(\beta\pi e). \quad (73)$$

By summing both terms, ELBO at \mathbf{x} can be estimated as

$$\begin{aligned} ELBO &= E_{\mathbf{x} \sim p(\mathbf{x})}[E_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(\cdot)] \\ &\simeq E_{\mathbf{x} \sim p(\mathbf{x})}[\log p(\mathbf{y})] \\ &\simeq E_{\mathbf{x} \sim p(\mathbf{x})}[\log p(\mathbf{x})]. \end{aligned} \quad (74)$$

As a result, ELBO (Eq. 1) in the original form (Kingma & Welling, 2014) is close to the log-likelihood of \mathbf{x} , regardless $\beta = 1$ or not, when the objective of β -VAE (Higgins et al., 2017) is optimised. Note that $\log p(\mathbf{x})$ in Eq. 73 is defined in the metric space. This also implies that the representation \mathbf{y} depends on the metrics.

Next, the predetermined conditional distribution $p_{\mathbb{R}^p}(\mathbf{x}|\hat{\mathbf{x}})$ used for training and the true conditional distribution $p_\theta(\mathbf{x}|\mathbf{z}) = p_\theta(\mathbf{x}|\hat{\mathbf{x}})$ after optimization are examined. Although $p_{\mathbb{R}^p}(\mathbf{x}|\hat{\mathbf{x}})$ and $p_\theta(\mathbf{x}|\hat{\mathbf{x}})$ are expected to be equivalent after optimization, the theoretical relationship between both is not well discussed. Assume $p_{\mathbb{R}^p}(\mathbf{x}|\hat{\mathbf{x}}) = \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \sigma^2 \mathbf{I})$. In this case, the metric $D(\mathbf{x}, \hat{\mathbf{x}})$ is derived as $-\log p_{\mathbb{R}^p}(\mathbf{x}|\hat{\mathbf{x}}) = (1/2\sigma^2)|\mathbf{x} - \hat{\mathbf{x}}|_2^2 + \text{Const}$. Using Eq. 53, the following equations are derived:

$$E_{p(\mathbf{x})}[D(\mathbf{x}, \hat{\mathbf{x}})] = E_{p(\mathbf{x})}[(1/2\sigma^2)|\mathbf{x} - \hat{\mathbf{x}}|_2^2] = E_{p(\mathbf{x})}[(1/2\sigma^2) \sum_i (x_i - \hat{x}_i)^2] \simeq n\beta/2. \quad (75)$$

Assume that $\sum_i (x_i - \hat{x}_i)^2$ for all i are equivalent. Then the next equation is derived:

$$E_{p(\mathbf{x})}[(x_i - \hat{x}_i)^2] \simeq \beta\sigma^2. \quad (76)$$

Because the variance of each dimension is $\beta\sigma^2$, the conditional distribution after optimization is estimated as $p_\theta(\mathbf{x}|\hat{\mathbf{x}}) = \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \beta\sigma^2 \mathbf{I})$.

If $\beta = 1$, i.e., the original VAE objective, both $p_{\mathbb{R}^p}(\mathbf{x}|\hat{\mathbf{x}})$ and $p_\theta(\mathbf{x}|\hat{\mathbf{x}})$ are equivalent. This result is consistent with what is expected.

If $\beta \neq 1$, however, $p_{\mathbb{R}^p}(\mathbf{x}|\hat{\mathbf{x}})$ and $p_\theta(\mathbf{x}|\hat{\mathbf{x}})$ are different. In other words, what β -VAE really does is to scale a variance of the pre-determined conditional distribution in the original VAE by a factor of β as Eq. 76. The detail is explained in Appendix B.3.

If $D(\mathbf{x}, \mathbf{x} + \delta\mathbf{x}) = {}^t\delta\mathbf{x}\mathbf{G}_\mathbf{x}\delta\mathbf{x} + O(|\delta\mathbf{x}|^3)$ is not SSE, by introducing a variable $\acute{\mathbf{x}} = \mathbf{L}_\mathbf{x}^{-1}\mathbf{x}$ where $\mathbf{L}_\mathbf{x}$ satisfies ${}^t\mathbf{L}_\mathbf{x}\mathbf{L}_\mathbf{x} = \mathbf{G}_\mathbf{x}$, the metric $D(\cdot, \cdot)$ can be replaced by SSE in the Euclidean space of $\acute{\mathbf{x}}$.

B.2. Relation to Tishby et al. (1999)

The theory described in Tishby et al. (1999), which first proposes the concept of information bottleneck (IB), is consistent with our analysis. Tishby et al. (1999) clarified the behaviour of the compressed representation when the rate-distortion trade-off is optimized. $\mathbf{x} \in X$ denotes the signal space with a fixed probability $p(\mathbf{x})$ and $\hat{\mathbf{x}} \in \hat{X}$ denotes its compressed representation. Let $D(\mathbf{x}, \hat{\mathbf{x}})$ be a loss metric. Then the rate-distortion trade-off can be described as:

$$L = I(X; \hat{X}) + \beta' E_{p(\mathbf{x}, \hat{\mathbf{x}})} [D(\mathbf{x}, \hat{\mathbf{x}})]. \quad (77)$$

By solving this condition, they derive the following equation:

$$p(\hat{\mathbf{x}}|\mathbf{x}) \propto \exp(-\beta' D(\mathbf{x}, \hat{\mathbf{x}})). \quad (78)$$

As shown in our discussion above, $p(\hat{\mathbf{x}}|\mathbf{x}) \simeq \mathcal{N}(\hat{\mathbf{x}}; \mathbf{x}, (\beta/2)\mathbf{I}_m)$ will hold in the metric defined space from our VAE analysis. This result is equivalent to Eq. 78 in their work if $D(\mathbf{x}, \hat{\mathbf{x}})$ is SSE and β' is set to β^{-1} , as follows:

$$p(\hat{\mathbf{x}}|\mathbf{x}) \propto \exp(-\beta' D(\mathbf{x}, \hat{\mathbf{x}})) = \exp\left(-\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2}{2(\beta/2)}\right) \propto \mathcal{N}(\hat{\mathbf{x}}; \mathbf{x}, (\beta/2)\mathbf{I}_m). \quad (79)$$

If $D(\mathbf{x}, \hat{\mathbf{x}})$ is not SSE, the use of the space transformation explained in appendix B.1 will lead to the same result.

B.3. Relation to β -VAE (Higgins et al., 2017)

This section explains the clear understanding of β -VAE (Higgins et al., 2017), and also corrects some of their theory.

In Higgins et al. (2017), ELBO equation is modified as:

$$E_{p(\mathbf{x})} [E_{\hat{\mathbf{x}} \sim p_\phi(\hat{\mathbf{x}}|\mathbf{x})} [q_\theta(\mathbf{x}|\hat{\mathbf{x}})] - \beta D_{\text{KL}}(\cdot)]. \quad (80)$$

However, they use the predetermined probabilities of $p_\theta(\hat{\mathbf{x}}|\mathbf{x})$ such as the Bernoulli and Gaussian distributions in training (described in table 1 in Higgins et al. (2017)). As shown in our appendix G.2, the log-likelihoods of the Bernoulli and Gaussian distributions can be regarded as BCE and SSE metrics, respectively. As a result, the actual objective for training in Higgins et al. (2017) is not Eq. 80, but the objective $L_x = D(\mathbf{x}, \hat{\mathbf{x}}) + \beta D_{\text{KL}}(\cdot)$ in Eq. 3 using BCE and SSE metrics with varying β . Thus ELBO as Eq. 1 form will become $\log p(\mathbf{x})$ in the BCE / SSE metric defined space regardless $\beta = 1$ or not, as shown in appendix B.1.

Actually, the equation 80 dose not show the log-likelihood of \mathbf{x} after optimization. When $D_{\text{KL}}(\cdot) \simeq -\log p(\mathbf{x}) - (n/2) \log(\beta\pi e)$ and $E_{\hat{\mathbf{x}} \sim p(\hat{\mathbf{x}}|\mathbf{x})} [p(\mathbf{x}|\hat{\mathbf{x}})] \simeq -(n/2) \log(\beta\pi e)$ are applied, the value of Eq. 80 is derived as $\beta \log p(\mathbf{x}) + (\beta - 1)(n/2) \log(\beta\pi e)$, which is different from the log-likelihood of \mathbf{x} in Eq. 73 if $\beta \neq 1$.

Correctly, what β -VAE really does is only to scale the variance of the pre-determined conditional distribution in the original VAE by a factor of β . In the case the pre-determined conditional distribution is Gaussian $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \sigma^2 \mathbf{I})$, the objective of β -VAE can be rewritten as a linearly scaled original VAE objective with a Gaussian $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \beta\sigma^2 \mathbf{I})$ where the variance is $\beta\sigma^2$ instead of σ^2 :

$$\begin{aligned} E_{q_\phi(\cdot)} [\log \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \sigma^2 \mathbf{I})] - \beta D_{\text{KL}}(\cdot) &= E_{q_\phi(\cdot)} \left[-\frac{1}{2} \log 2\pi\sigma^2 - \frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2}{2\sigma^2} \right] - \beta D_{\text{KL}}(\cdot) \\ &= \beta \left(E_{q_\phi(\cdot)} \left[-\frac{1}{2} \log 2\pi\beta\sigma^2 - \frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2}{2\beta\sigma^2} \right] - D_{\text{KL}}(\cdot) \right) \\ &\quad + \frac{\beta}{2} \log 2\pi\beta\sigma^2 - \frac{1}{2} \log 2\pi\sigma^2 \\ &= \beta \left(\underline{E_{q_\phi(\cdot)} [\log \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \beta\sigma^2 \mathbf{I})]} - D_{\text{KL}}(\cdot) \right) + \text{const.} \end{aligned} \quad (81)$$

Here, the underlined terms in the last equation is just the ELBO with the predetermined conditional distribution $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \beta\sigma^2 \mathbf{I})$. So the optimization of β -VAE objective with the predetermined conditional distribution $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \sigma^2 \mathbf{I})$ is just the same as the optimization of the original VAE objective ($\beta=1$) with the predetermined conditional distribution $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \beta\sigma^2 \mathbf{I})$.

B.4. Relation to Alemi et al. (2018)

Alemi et al. (2018) discuss the rate-distortion trade-off by the theoretical entropy analysis. Their work is also presumed that the objective L_x was not mistakenly distinguished from ELBO, which leads to the incorrect discussion. In their work, the differential entropy for the input H , distortion D , and rate R are derived carefully. They suggest that VAE with $\beta = 1$ is sensitive (unstable) because D and R can be arbitrary value on the line $R = H - \beta D = H - D$. Furthermore, they also suggest that $R \geq H$, $D = 0$ at $\beta \rightarrow 0$ and $R = 0$, $D \geq H$ at $\beta \rightarrow \infty$ will hold as shown the figure 1 of their work.

In this appendix, we will show that β determines the value of R and D specifically. We also show that $R \simeq H - D$ will hold regardless $\beta = 1$ or not.

In their work, these values of H , D , and D are mathematically defined as:

$$H \equiv - \int d\mathbf{x} p^*(\mathbf{x}) \log p^*(\mathbf{x}), \quad (82)$$

$$D \equiv - \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{z} e(\mathbf{z}|\mathbf{x}) \log d(\mathbf{x}|\mathbf{z}), \quad (83)$$

$$R \equiv \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{z} e(\mathbf{z}|\mathbf{x}) \log \frac{e(\mathbf{z}|\mathbf{x})}{m(\mathbf{z})}. \quad (84)$$

Here, $p^*(\mathbf{x})$ is a true PDF of \mathbf{x} , $e(\mathbf{z}|\mathbf{x})$ is a stochastic encoder, $e(\mathbf{z}|\mathbf{x})$ is a decoder, and $m(\mathbf{z})$ is a marginal probability of \mathbf{z} .

Our work allows a rough estimation of Eqs. 82-84 with β by introducing the implicit isometric variable \mathbf{y} as explained in our work.

Using isometric variable \mathbf{y} and the relation $d\mathbf{z} e(\mathbf{z}|\mathbf{x}) = d\mathbf{y} e(\mathbf{y}|\mathbf{x})$, Eq. 83 can be rewritten as:

$$D = - \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{y} e(\mathbf{y}|\mathbf{x}) \log d(\mathbf{x}|\mathbf{y}). \quad (85)$$

Let μ_y be the implicit isometric latent variable corresponding to the mean of encoder output $\mu_{(x)}$. As discussed in section 4.1, $e(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mu_y, (\beta/2)\mathbf{I}_n)$ will hold. Because of isometricity, the value of $d(\mathbf{x}|\mathbf{y})$ will be also close to $e(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mu_y, (\beta/2)\mathbf{I}_n)$. Though $d(\mathbf{x}|\mathbf{z})$ must depend on $e(\mathbf{z}|\mathbf{x})$, this important point has not been discussed well in this work. By using the implicit isometric variable, we can connect both theoretically. Thus, D can be estimated as:

$$\begin{aligned} D &\simeq \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{y} \mathcal{N}(\mathbf{y}; \mu_y, (\beta/2)\mathbf{I}_n) \log \mathcal{N}(\mathbf{y}; \mu_y, (\beta/2)\mathbf{I}_n) \\ &\simeq \int d\mathbf{x} p^*(\mathbf{x}) \left(\frac{n}{2} \log(\beta\pi e) \right) \\ &= \frac{n}{2} \log(\beta\pi e). \end{aligned} \quad (86)$$

Second, R is examined. $m(\mathbf{y})$ is a marginal probability of \mathbf{y} . Using the relation $d\mathbf{z} e(\mathbf{z}|\mathbf{x}) = d\mathbf{y} e(\mathbf{y}|\mathbf{x})$ and $e(\mathbf{z}|\mathbf{x})/m(\mathbf{z}) = (e(\mathbf{y}|\mathbf{x})(d\mathbf{y}/d\mathbf{z}))/m(\mathbf{y})(d\mathbf{y}/d\mathbf{z}) = e(\mathbf{y}|\mathbf{x})/m(\mathbf{y})$, Eq. 84 can be rewritten as:

$$R \simeq \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{y} e(\mathbf{y}|\mathbf{x}) \log \frac{e(\mathbf{y}|\mathbf{x})}{m(\mathbf{y})}. \quad (87)$$

Because of isometricity, $e(\mathbf{y}|\mathbf{x}) \simeq p(\hat{\mathbf{x}}|\mathbf{x}) \simeq \mathcal{N}(\hat{\mathbf{x}}; \mathbf{x}, (\beta/2)\mathbf{I}_m)$ will approximately hold where $\hat{\mathbf{x}}$ denotes a decoder output. Thus $m(\mathbf{y})$ can be approximated by:

$$m(\mathbf{y}) \simeq \int d\mathbf{x} p^*(\mathbf{x}) e(\mathbf{y}|\mathbf{x}) \simeq \int d\mathbf{x} p^*(\mathbf{x}) \mathcal{N}(\hat{\mathbf{x}}; \mathbf{x}, (\beta/2)\mathbf{I}_m) \quad (88)$$

Here, if $\beta/2$, i.e., added noise, is small enough compared to the variance of \mathbf{x} , a normal distribution function term in this equation will act like a delta function. Thus $m(\mathbf{y})$ can be approximated as:

$$m(\mathbf{y}) \simeq \int d\hat{\mathbf{x}} p^*(\hat{\mathbf{x}}) \delta(\hat{\mathbf{x}} - \mathbf{x}) \simeq p^*(\mathbf{x}). \quad (89)$$

In the similar way, the following approximation will also hold.

$$\int d\mathbf{y} e(\mathbf{y}|\mathbf{x}) \log m(\mathbf{y}) \simeq \int d\mathbf{y} e(\mathbf{y}|\mathbf{x}) \log p^*(\mathbf{x}) \simeq \int d\mathbf{x} \delta(\mathbf{x} - \mathbf{x}) \log p^*(\mathbf{x}) \simeq \log p^*(\mathbf{x}). \quad (90)$$

By using these approximation and applying Eqs. 85-86, R in Eq. 84 can be approximated as:

$$\begin{aligned} R &\simeq \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{y} e(\mathbf{y}|\mathbf{x}) \log \frac{e(\mathbf{y}|\mathbf{x})}{p^*(\mathbf{x})} \\ &\simeq - \int d\mathbf{x} p^*(\mathbf{x}) \log p^*(\mathbf{x}) - \left(- \int d\mathbf{x} p^*(\mathbf{x}) \int d\mathbf{y} e(\mathbf{y}|\mathbf{x}) \log e(\mathbf{y}|\mathbf{x}) \right) \\ &\simeq H - \frac{n}{2} \log(\beta \pi e) \\ &\simeq H - D. \end{aligned} \quad (91)$$

As discussed above, R and D can be specifically derived from β . In addition, Shannon lower bound discussed in Alemi et al. (2018) can be roughly verified in the optimized VAE with clearer notations using β .

From the discussion above, we presume Alemi et al. (2018) might wrongly treat D in their work. They suggest that VAE with $\beta = 1$ is sensitive (unstable) because D and R can be arbitrary value on the line $R = H - \beta D = H - D$; however, our work as well as Tishby et al. (1999) (appendix B.2) and Dai & Wipf (2019) (appendix B.5) show that the differential entropy of the distortion and rate, i.e., D and R , are specifically determined by β after optimization, and $R = H - D$ will hold for any β regardless $\beta = 1$ or not. Alemi et al. (2018) also suggest D should satisfy $D \geq 0$ because D is a distortion; however, we suggest D should be treated as a differential entropy and can be less than 0 because \mathbf{x} is once handled as a continuous signal with a stochastic process in Eqs. 82-84. Here, $D \simeq (n/2) \log(\beta \pi e)$ can be $-\infty$ if $\beta \rightarrow 0$, as also shown in Dai & Wipf (2019). Thus, upper bound of R at $\beta \rightarrow 0$ is not H , but $R = H - (-\infty) = \infty$, as shown in RD theory for a continuous signal. Huang et al. (2020) show this property experimentally in their figures 4-8 such that R seems to diverge if MSE is close to 0.

B.5. Relation to Dai et al. (2018) and Dai & Wipf (2019)

Our work is consistent with Dai et al. (2018) and Dai & Wipf (2019).

Dai et al. (2018) analyses VAE by assuming a linear model. As a result, the estimated posterior is constant. If the distribution of the manifold is the Gaussian, our work and Dai et al. (2018) give a similar result with constant posterior variances. For non-Gaussian data, however, the quantitative analysis such as probability estimation is intractable using their linear model. Our work reveals that the posterior variance gives a scaling factor between \mathbf{z} in VAE and \mathbf{y} in the isometric space when VAE is ideally trained with rich parameters. This is validated by Figures 3c and 3d, where the estimation of the posterior variance at each data point is a key.

Next, the relation to Dai & Wipf (2019) is discussed. They analyse a behavior of VAE when ideally trained. For example, the theorem 5 in their work shows that $D \rightarrow (d/2) \log \gamma + O(1)$ and $R \rightarrow -(\hat{\gamma}/2) \log \gamma + O(1)$ hold if $\gamma \rightarrow +0$, where γ , d , and $\hat{\gamma}$ denote a variance of $d(\mathbf{x}|\mathbf{z})$, data dimension, and latent dimension, respectively. By setting $\gamma = \beta/2$ and $d = \hat{\gamma} = n$, this theorem is consistent with R and D derived in Eq. 86 and Eq. 91.

B.6. Relation to Rate-distortion theory (Berger, 1971) and transform coding (Goyal, 2001; Pearlman & Said, 2011)

RD theory (Berger, 1971) formulated the optimal transform coding (Goyal, 2001; Pearlman & Said, 2011) for the Gaussian source with square error metric as follows. Let $\mathbf{x} \in \mathbb{R}^m$ be a point in a dataset. First, the data are transformed deterministically with the orthonormal transform (orthogonal and unit norm) such as Karhunen-Loève transform (KLT) (Rao & Yip, 2000). Note that the basis of KLT is equivalent to a PCA basis. Let $\mathbf{z} \in \mathbb{R}^m$ be a point transformed from \mathbf{x} . Then, \mathbf{z} is entropy-coded by allowing equivalent stochastic distortion (or posterior with constant variance) in each dimension. A lower bound of a rate R at a distortion D is denoted by $R(D)$. The derivation of $R(D)$ is as follows. Let z_j be the j -th dimensional component of \mathbf{z} and $\sigma_{z_j}^2$ be the variance of z_j in a dataset. It is noted that $\sigma_{z_j}^2$ is the equivalent to eigenvalues of PCA for the dataset. Let d be a distortion equally allowed in each dimensional channel. At the optimal condition, the

distortion D_{opt} and rate R_{opt} on the curve $R(D)$ is calculated as a function of d :

$$\begin{aligned} R_{\text{opt}} &= \frac{1}{2} \sum_{j=1}^m \max(\log(\sigma_{z_j}^2/d), 0), \\ D_{\text{opt}} &= \sum_{j=1}^m \min(d, \sigma_{z_j}^2). \end{aligned} \quad (92)$$

The simplest way to allow equivalent distortion is to use a uniform quantization (Goyal, 2001). Let T be a quantization step, and $\text{round}(\cdot)$ be a round function. Quantized value \hat{z}_j is derived as kT , where $k = \text{round}(z_j/T)$. Then, d is approximated by $T^2/12$ as explained in Appendix G.1.

To practically achieve the best RD trade-off in image compression, rate-distortion optimization (RDO) has also been widely used (Sullivan & Wiegand, 1998). In RDO, the best trade-off is achieved by finding a encoding parameter that minimizes the cost L at given Lagrange parameter λ as:

$$L = D + \lambda R. \quad (93)$$

This equation is equivalent to VAE when $\lambda = \beta^{-1}$.

We show the optimum condition of VAE shown in Eq. 62 and 65 can be mapped to the optimum condition of transform coding (Goyal, 2001) as shown in Eq. 92. First, the derivation of Eq. 92 is explained by solving the optimal distortion assignment to each dimension. In the transform coding for m dimensional the Gaussian data, an input data \mathbf{x} is transformed to \mathbf{z} using an orthonormal transform such as KLT/DCT. Then each dimensional component z_j is encoded with allowing distortion d_j . Let D be a target distortion satisfying $D = \sum_{j=1}^m d_j$. Next, $\sigma_{z_j}^2$ denotes a variance of each dimensional component z_j for the input dataset. Then, a rate R can be derived as $\sum_{j=1}^m \frac{1}{2} \log(\sigma_{z_j}^2/d_j)$. By introducing a Lagrange parameter λ and minimizing a rate-distortion optimization cost $L = D + \lambda R$, the optimum condition is derived as:

$$\lambda_{\text{opt}} = 2D/m, \quad d_j = D/m = \lambda_{\text{opt}}/2. \quad (94)$$

This result is consistent with Eq. 62 and 65 by setting $\beta = \lambda_{\text{opt}} = 2D/m$. This implies that $L_G = D_G + \beta R_G$ is a rate-distortion optimization (RDO) cost of transform coding when \mathbf{x} is deterministically transformed to \mathbf{y} in the implicit isometric space and stochastically encoded with a distortion $\beta/2$.

C. Details of the networks and training conditions in the experiments

This appendix explains the networks and training conditions in Section 5.

C.1. Toy data set

This appendix explains the details of the networks and training conditions in the experiment of the toy data set in Section 5.1.

Network configurations:

FC(i, o, f) denotes a FC layer with input dimension i, output dimension o, and activate function f.

The encoder network is composed of FC(16, 128, tanh)-FC(128, 64, tanh)-FC(64, 3, linear) $\times 2$ (for μ and σ). The decoder network is composed of FC(3, 64, tanh)-FC(64, 128, tanh)-FC(128, 16, linear).

Training conditions:

The reconstruction loss $D(\cdot, \cdot)$ is derived such that the loss per input dimension is calculated and all of the losses are averaged by the input dimension $m = 16$. The KL divergence is derived as a summation of $D_{\text{KL}(j)}(\cdot)$ as explained in Eq. 2.

In our code, we use essentially the same, but a constant factor scaled loss objective from the original β -VAE form $L_{\mathbf{x}} = D(\cdot, \cdot) + \beta D_{\text{KL}(j)}(\cdot)$ in Eq. 1, such as:

$$L_{\mathbf{x}} = \lambda D(\cdot, \cdot) + D_{\text{KL}(j)}(\cdot). \quad (95)$$

Equation 95 is essentially equivalent to $L = D(\cdot, \cdot) + \beta D_{\text{KL}(j)}(\cdot)$, multiplying a constant $\lambda = \beta^{-1}$ to the original form. The reason why we use this form is as follows. Let $\text{ELBO}_{\text{true}}$ be the true ELBO in the sense of log-likelihood, such as $E[\log p(\mathbf{x})]$. As shown in Eq. 57, the minimum of the loss objective in the original β -VAE form is likely to be a

$-\beta \text{ELBO}_{\text{true}} + \text{Constant}$. If we use Eq. 95, the minimum of the loss objective will be $-\text{ELBO}_{\text{true}} + \text{Constant}$, which seems more natural form of ELBO. Thus, Eq. 95 allows estimating a data probability from $L_{\mathbf{x}}$ in Eqs. 19 and 21, without scaling $L_{\mathbf{x}}$ by $1/\beta$.

Then the network is trained with $\lambda = \beta^{-1} = 100$ using 500 epochs with a batch size of 128. Here, Adam optimizer is used with the learning rate of 1e-3. We use a PC with CPU Inter(R) Xeon(R) CPU E3-1280v5@3.70GHz, 32GB memory equipped with NVIDIA GeForce GTX 1080. The simulation time for each trial is about 20 minutes, including the statistics evaluation codes.

C.2. CelebA data set

This appendix explains the details of the networks and training conditions in the experiment of the toy data set in Section 5.2.

Network configurations:

CNN(w, h, s, c, f) denotes a CNN layer with kernel size (w, h), stride size s, dimension c, and activate function f. GDN and IGDN[†] are activation functions designed for image compression (Ballé et al., 2016). This activation function is effective and popular in deep image compression studies.

The encoder network is composed of CNN(9, 9, 2, 64, GDN) - CNN(5, 5, 2, 64, GDN) - CNN(5, 5, 2, 64, GDN) - CNN(5, 5, 2, 64, GDN) - FC(1024, 1024, softplus) - FC(1024, 32, None) $\times 2$ (for μ and σ) in encoder.

The decoder network is composed of FC(32, 1024, softplus) - FC(1024, 1024, softplus) - CNN(5, 5, 2, 64, IGDN) - CNN(5, 5, 2, 64, IGDN) - CNN(5, 5, 2, 64, IGDN)-CNN(9, 9, 2, 3, IGDN).

Training conditions:

In this experiment, SSIM explained in Appendix G.2 is used as a reconstruction loss. The reconstruction loss $D(\cdot, \cdot)$ is derived as follows. Let SSIM be a SSIM calculated from two input images. As explained in Appendix G.2, SSIM is measured for a whole image, and its range is between 0 and 1. If the quality is high, SSIM value becomes close to 1. Then $1 - \text{SSIM}$ is set to $D(\cdot, \cdot)$.

We also use the loss form as in Equation 95 in our code. In the case of the decomposed loss, the loss function $L_{\mathbf{x}}$ is set to $\lambda(D(\mathbf{x}, \hat{\mathbf{x}}) + D(\hat{\mathbf{x}}, \mathbf{x})) + D_{\text{KL}}(\cdot)$ in our code. Then, the network is trained with $\lambda = \beta^{-1} = 1,000$ using a batch size of 64 for 300,000 iterations. Here, Adam optimizer is used with the learning rate of 1e-3.

We use a PC with CPU Intel(R) Core(TM) i7-6850K CPU @ 3.60GHz, 12GB memory equipped with NVIDIA GeForce GTX 1080. The simulation time for each trial is about 180 minutes, including the statistics evaluation codes.

[†]Google provides a code in the official Tensorflow library (<https://github.com/tensorflow/compression>)

D. Additional results in the toy datasets

D.1. Scattering plots for the square error loss in Section

Figure 10a shows the plots of $p(\mathbf{x})$ and estimated probabilities for the square error coding loss in Section 5.1, where the scale factor $a_{\mathbf{x}}$ in Eq. 21 is 1. Thus, both $\exp(-L_{\mathbf{x}}/\beta)$ and $p(\boldsymbol{\mu}_{(\mathbf{x})}) \prod_j \sigma_{j(\mathbf{x})}$ show a high correlation, allowing easy estimation of the data probability in the input space. In contrast, $p(\boldsymbol{\mu}_{(\mathbf{x})})$ still shows a low correlation. These results are consistent with our theory.

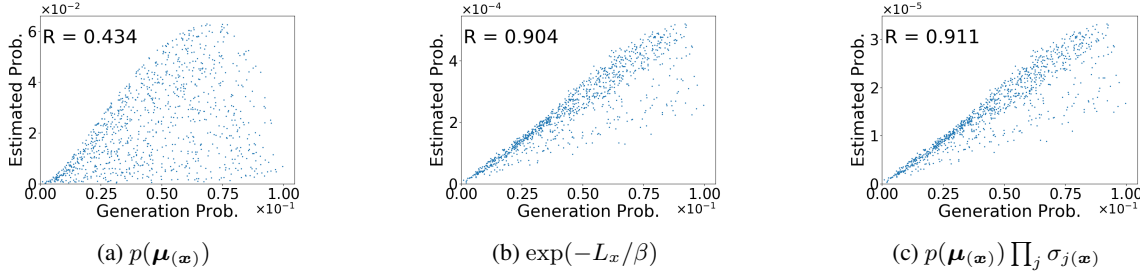


Figure 10. Plots of the data generation probability (x-axis) versus estimated probabilities (y-axes) for the square error loss. y-axes are (a) $p(\boldsymbol{\mu}_{(\mathbf{x})})$, (b) $\exp(-L_{\mathbf{x}}/\beta)$, and (c) $p(\boldsymbol{\mu}_{(\mathbf{x})}) \prod_j \sigma_{j(\mathbf{x})}$.

D.2. Ablation study using 3 toy datasets, 3 coding losses, and 10 β parameters.

In this appendix, we explain the ablation study for the toy datasets. We introduce three toy datasets and three coding losses including those used in Section 5.1. We also change $\beta^{-1} = \lambda$ from 1 to 1,000 in training. The details of the experimental conditions are shown as follows.

Datasets: First, we call the toy dataset used in Section 5.1 the Mix dataset in order to distinguish three datasets. The second dataset is generated such that three dimensional variables s_1 , s_2 , and s_3 are sampled in accordance with the distributions $p(s_1)$, $p(s_2)$, and $p(s_3)$ in Figure 11. The variances of the variables are the same as those of the Mix dataset, i.e., $1/6$, $2/3$, and $8/3$, respectively. We call this the Ramp dataset. Because the PDF shape of this dataset is quite different from the prior $\mathcal{N}(\mathbf{z}; 0, I_3)$, the fitting will be the most difficult among the three. The third dataset is generated such that three dimensional variables s_1 , s_2 , and s_3 are sampled in accordance with the normal distributions $\mathcal{N}(s_1; 0, 1/6)$, $\mathcal{N}(s_2; 0, 2/3)$, and $\mathcal{N}(s_3; 0, 8/3)$, respectively. We call this the Norm dataset. The fitting will be the easiest, because both the prior and input have the normal distributions, and the posterior standard deviation, given by the PDF ratio at the same CDF, can be a constant.

Coding losses: Two of the three coding losses is the square error loss and the downward-convex loss described in Section 5.1. The third coding loss is an upward-convex loss which we design as Eq. 96 such that the scale factor $a_{\mathbf{x}}$ becomes the reciprocal of the scale factor in Eq. 24:

$$D(\mathbf{x}, \hat{\mathbf{x}}) = a_{\mathbf{x}} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2, \quad \text{where } a_{\mathbf{x}} = (2/3 + 2 \|\mathbf{x}\|_2^2/21)^{-1} \text{ and } \mathbf{G}_{\mathbf{x}} = a_{\mathbf{x}} \mathbf{I}_m. \quad (96)$$

Figure 12 shows the scale factors $a_{\mathbf{x}}$ in Eqs. 24 and 96, where s_1 in $\mathbf{x} = (s_1, 0, 0)$ moves within ± 5 .

Parameters: As explained in Appendix C.1, $\lambda = 1/\beta$ is used as a hyper parameter. Specifically, $\lambda = 1, 2, 5, 10, 20, 50, 100, 200, 500$, and 1,000 are used.

Figures 13 - 21 show the property measurements for all combinations of the datasets and coding losses, with changing λ . In each Figure, the estimated norms of the implicit transform are shown in the figure (a), the ratios of the estimated variances are shown in the figure (b), and the correlation coefficients between $p(\mathbf{x})$ and estimated data probabilities are shown in the figure (c), respectively.

First, the estimated norm of the implicit transform in the figures (a) is discussed. In all conditions, the norms are close to 1 as described in Eq. 23 in the λ range 50 to 1000. These results show consistency with our theoretical analysis, supporting the existence of the implicit orthonormal transform. The values in the Norm dataset are the closest to 1, and those in the Ramp dataset are the most different, which seems consistent with the difficulty of the fitting.

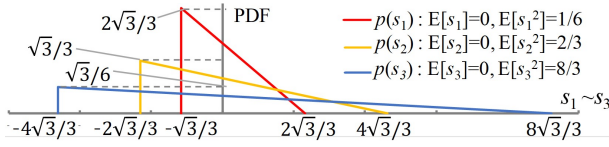
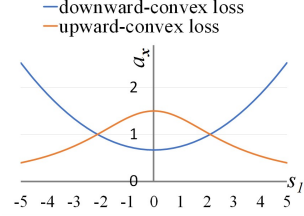


Figure 11. PDFs of three variables to generate a Ramp dataset.


 Figure 12. Scale factor a_x for the downward-convex loss and upward-convex loss.

Second, the ratio of the estimated variances is discussed. In the figures (b), $\text{Var}(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$. Then, $\text{Var}(z_2)/\text{Var}(z_1)$ and $\text{Var}(z_3)/\text{Var}(z_1)$ are plotted. In all conditions, the ratios of $\text{Var}(z_2)/\text{Var}(z_1)$ and $\text{Var}(z_3)/\text{Var}(z_1)$ are close to the variance ratios of the input variables, i.e., 4 and 16, in the λ range 5 to 500. Figure 22 shows the detailed comparison of the ratio for the three datasets and three coding losses at $\lambda = 100$. In most cases, the estimated variances in the downward-convex loss are the smallest, and those in the upward-convex loss are the largest, which is more distinct for $\text{Var}(z_3)/\text{Var}(z_1)$. This can be explained as follows. When using the downward-convex loss, the space region with a large norm is thought of as shrinking in the inner product space, as described in Section 5.1. This will make the variance smaller. In contrast, when using the upward-convex loss, the space region with a large norm is thought of as expanding in the inner product space, making the variance larger. Here, the dependency of the losses on the ratio changes is less in the Norm dataset. The possible reason is that data in the normal distribution concentrate around the center, having less effect on the loss scale factor in the downward-convex loss and upward-convex loss.

Third, the correlation coefficients between $p(x)$ and the estimated data probabilities in the figures (c) are discussed. In the Mix dataset and Ramp dataset, the correlation coefficients are around 0.9 in the λ range from 20 to 200 when the estimated probabilities $a_x^{n/2} p(\mu(x)) \prod_{j=1}^n \sigma_{j(x)}$ and $a_x^{n/2} \exp(-(1/\beta)L_x)$ in Eq. 21 are used. When using $p(\mu(x)) \prod_{j=1}^n \sigma_{j(x)}$ and $\exp(-(1/\beta)L_x)$ in the downward-convex loss and upward-convex loss, the correlation coefficients become worse. In addition, when using the prior probability $p(\mu(x))$, the correlation coefficients always show the worst. In the Norm dataset, the correlation coefficients are close to 1.0 in the wider range of λ when using the estimated distribution in Eq. 21. When using $p(\mu(x)) \prod_{j=1}^n \sigma_{j(x)}$ and $\exp(-(1/\beta)L_x)$ in the downward-convex loss and upward-convex loss, the correlation coefficients also become worse. When using the prior probability $p(\mu(x))$, however, the correlation coefficients are close to 1 in contrast to the other two datasets. This can be explained because both the input distribution and the prior distribution are the same normal distribution, allowing the posterior variances almost constant. These results also show consistency with our theoretical analysis.

Figure 23 shows the dependency of the coding loss on β for the Mix, Ramp, and Norm dataset using square the error loss. From D_G in Eq. 64 and $n = 3$, the theoretical value of coding loss is $\frac{3\beta}{2}$, as also shown in the figure. Unlike Figs. 13-21, x -axis is $\beta = \lambda^{-1}$ to evaluate the linearity. As expected in Theorem 3, the coding losses are close to the theoretical value where $\beta < 0.1$, i.e., $\lambda > 10$.

Figure 24 shows the dependency of the ratio of transform loss to coding loss on β for the Mix, Ramp, and Norm dataset using square the error loss. From Eq. 43, the estimated transform loss is $\sum_{i=1}^3 (\beta/2)^2 / \text{Var}(s_i) = \frac{63\beta^2}{32}$. Thus the theoretical value is $(\frac{63\beta^2}{32}) / (\frac{3\beta}{2}) = \frac{21\beta}{16}$, as is also shown in the figure. x -axis is also $\beta = \lambda^{-1}$ like Figure 23. Considering the correlation coefficient discussed above, the useful range of β seems between 0.005-0.05 (20-200 for λ). In this range, the ratio is less than 0.1, implying the transform loss is almost negligible. As expected in Lemma 3 and appendix A.3, the ratio is close to the theoretical value where $\beta > 0.01$, i.e., $\lambda < 100$. For $\beta < 0.01$, the transform loss is still negligibly small, but the ratio is somewhat off the theoretical value. The reason is presumably that the transform loss is too small to fit the network.

As shown above, this ablation study strongly supports our theoretical analysis in sections 4.

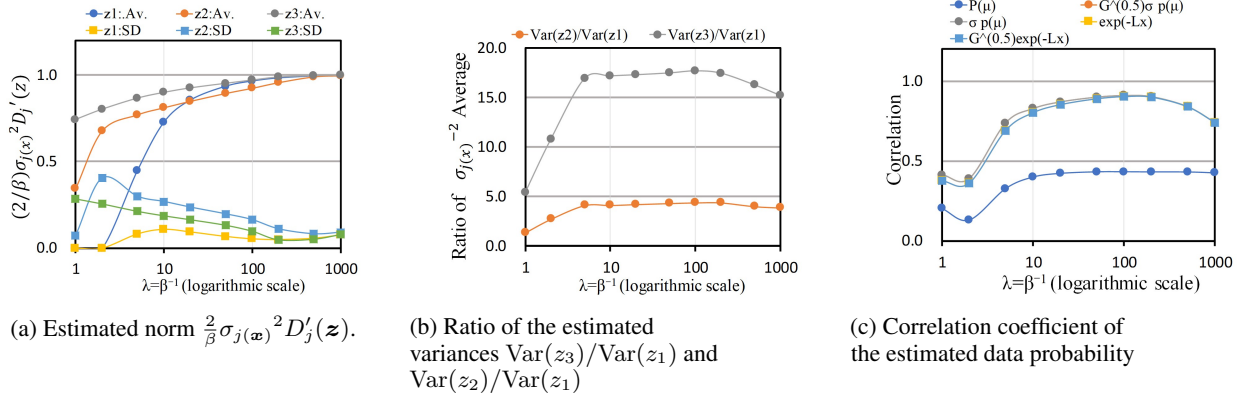


Figure 13. Property measurements of the Mix dataset using the square error loss. λ is changed from 1 to 1,000. $Var(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$.

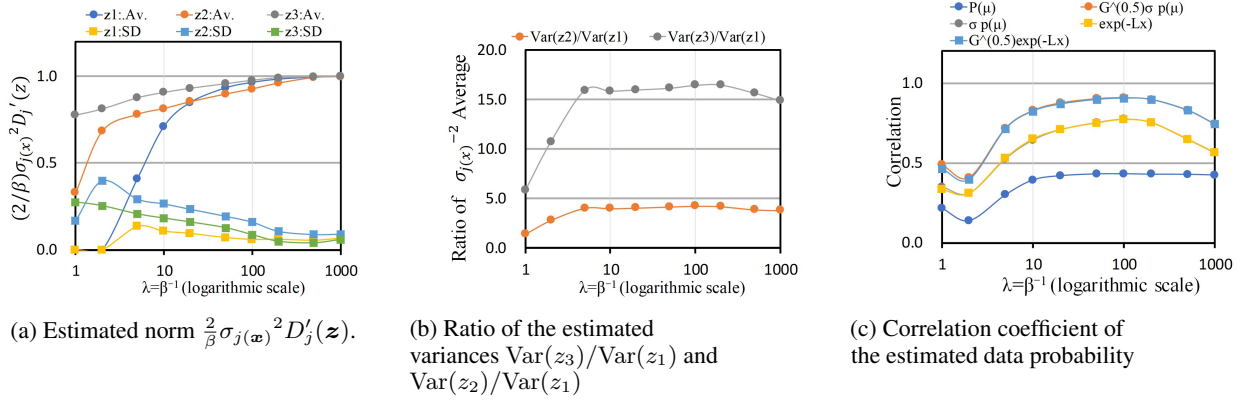


Figure 14. Property measurements of the Mix dataset using the downward-convex loss. λ is changed from 1 to 1,000. $Var(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$.

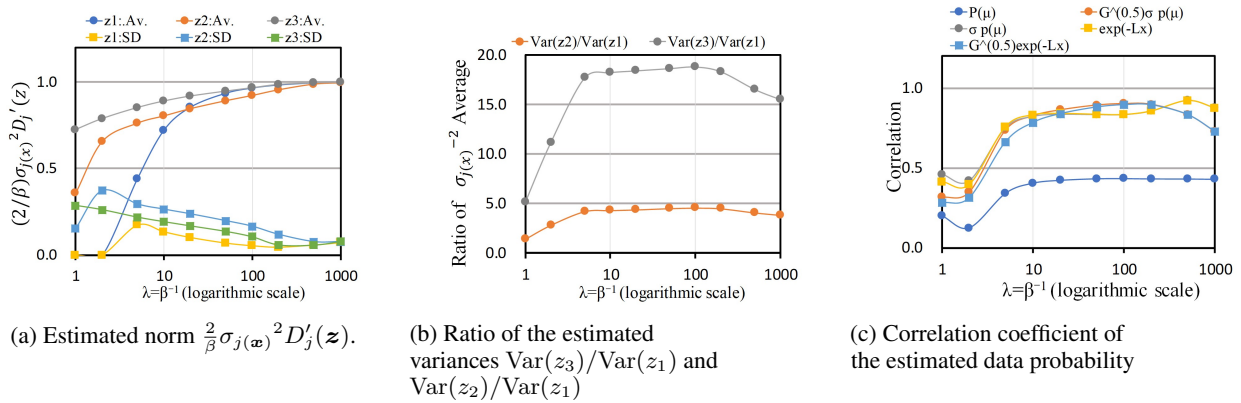


Figure 15. Property measurements of the Mix dataset using the upward-convex loss. λ is changed from 1 to 1,000. $Var(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$.

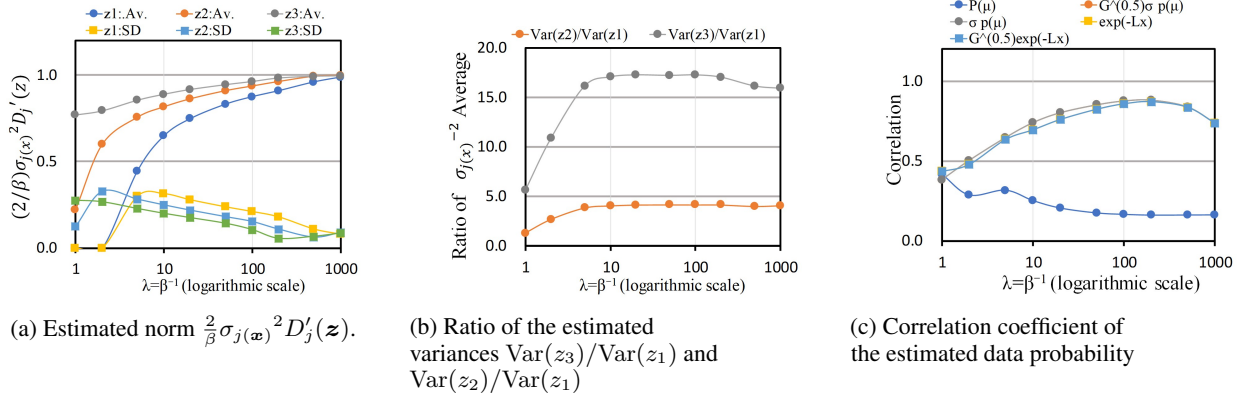


Figure 16. Property measurements of the Ramp dataset using the square error loss. λ is changed from 1 to 1,000. $\text{Var}(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(\mathbf{x})}^{-2}$.

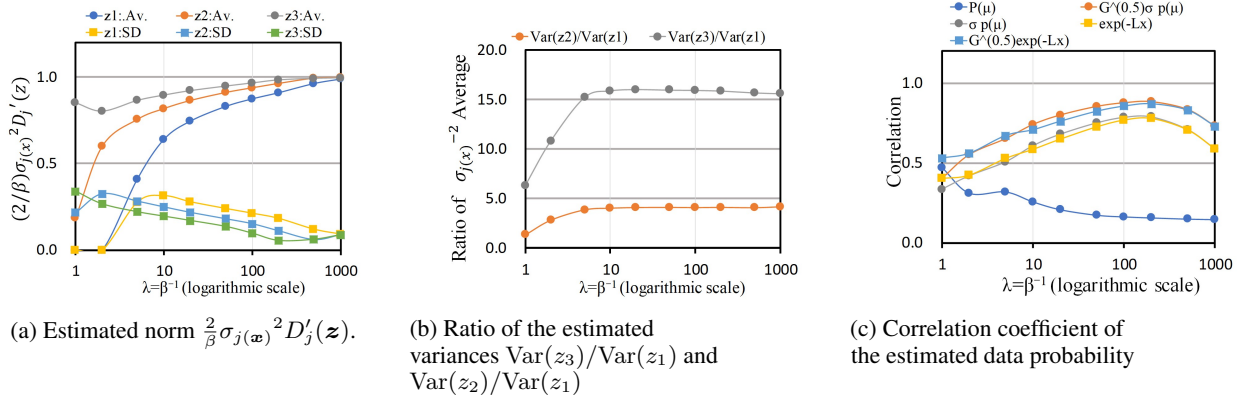


Figure 17. Property measurements of the Ramp dataset using the downward-convex loss. λ is changed from 1 to 1,000. $\text{Var}(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(\mathbf{x})}^{-2}$.

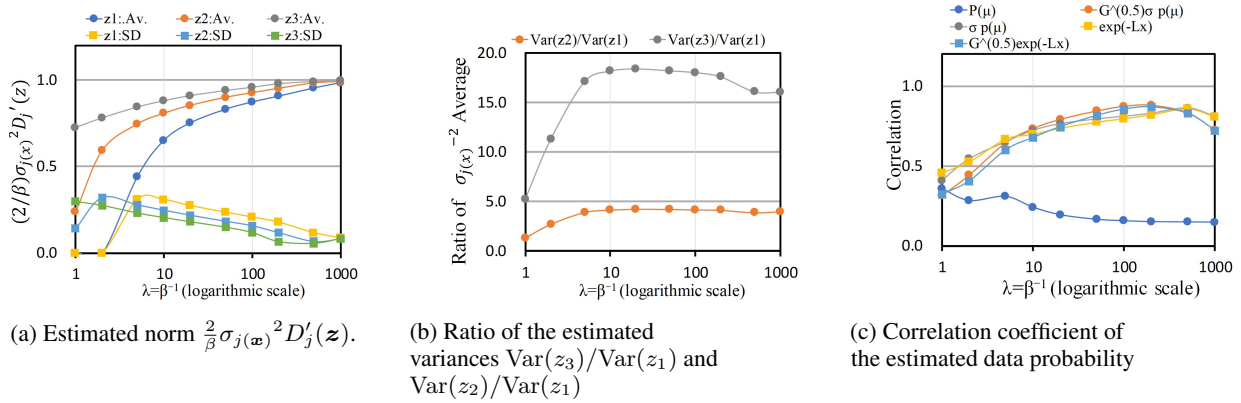


Figure 18. Property measurements of the Ramp dataset using the upward-convex loss. λ is changed from 1 to 1,000. $\text{Var}(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(\mathbf{x})}^{-2}$.

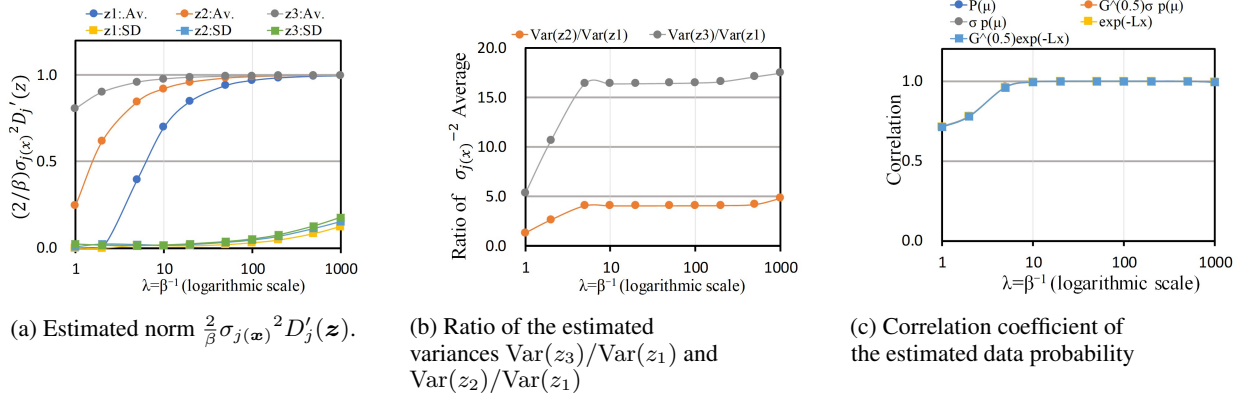


Figure 19. Property measurements of the Norm dataset using the square error loss. λ is changed from 1 to 1,000. $Var(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$.

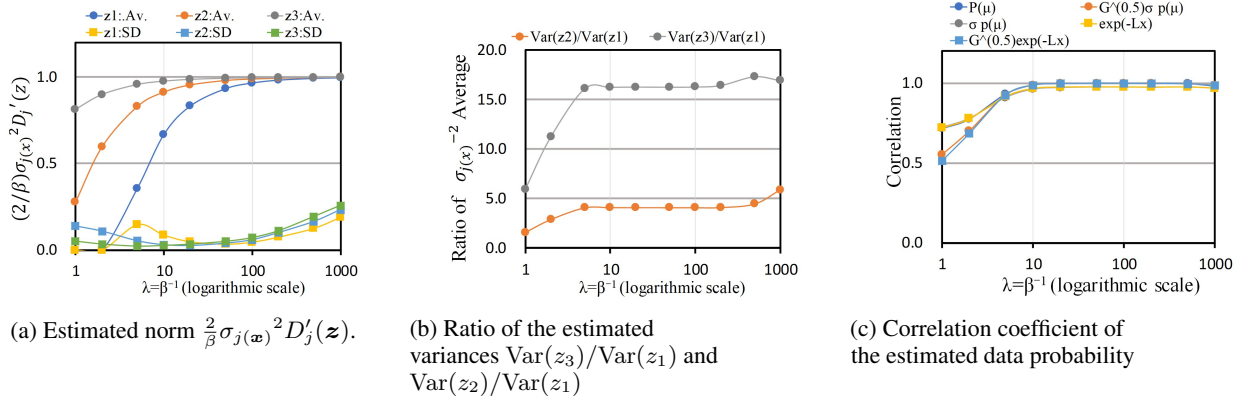


Figure 20. Property measurements of the Norm dataset using the downward-convex loss. λ is changed from 1 to 1,000. $Var(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$.

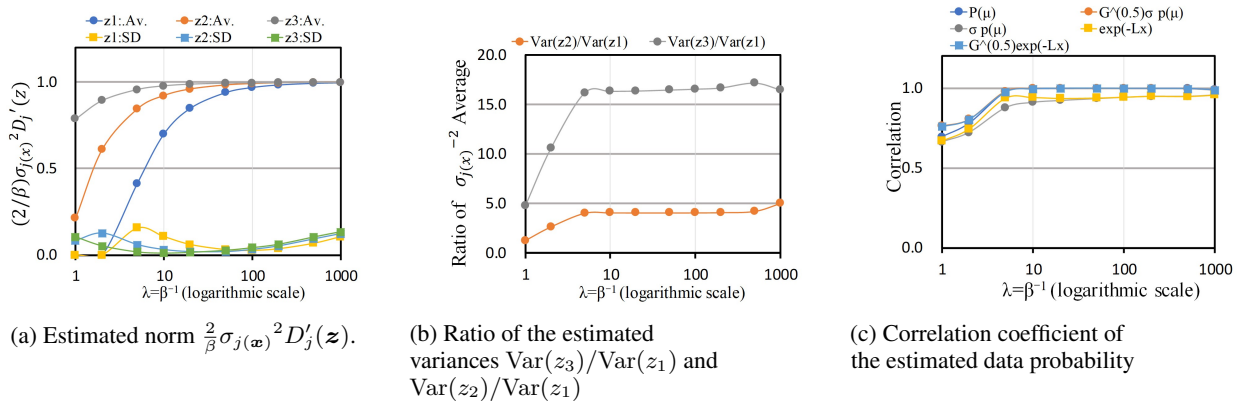


Figure 21. Property measurements of the Mix dataset using the upward-convex loss. λ is changed from 1 to 1,000. $Var(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(x)}^{-2}$.

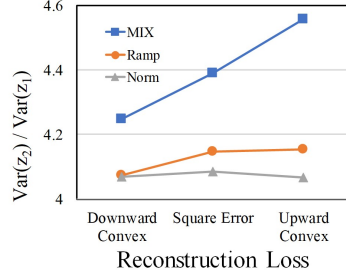
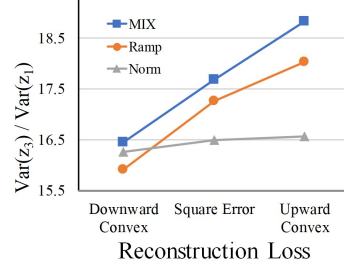
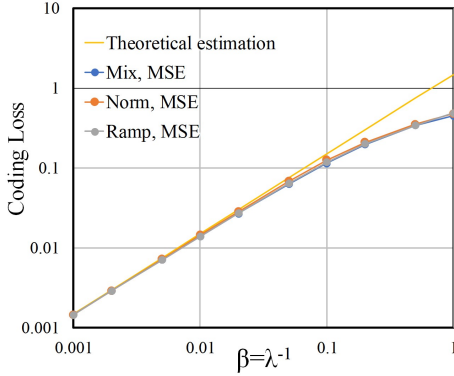
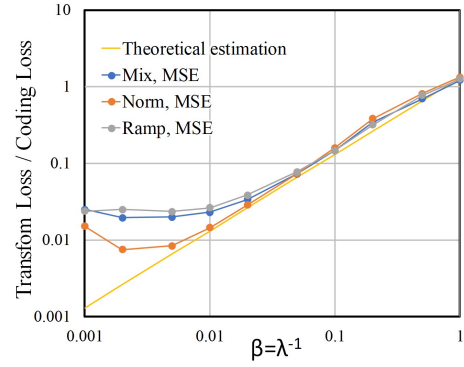

 (a) $\text{Var}(z_2)/\text{Var}(z_1)$.

 (b) $\text{Var}(z_3)/\text{Var}(z_1)$.

 Figure 22. Ratio of the estimated variances $\text{Var}(z_3)/\text{Var}(z_1)$ and $\text{Var}(z_2)/\text{Var}(z_1)$ for the three datasets and three coding losses at $\lambda = 100$. $\text{Var}(z_j)$ denotes the estimated variance, given by the average of $\sigma_{j(\mathbf{x})}^{-2}$.

 Figure 23. Dependency of Coding Loss on β for Mix, Norm, and Ramp dataset using square loss.

 Figure 24. Dependency of Transform loss / Coding Loss Ratio on β for Mix, Norm, and Ramp dataset using square loss.

D.3. Increase of latent dimension

The Table 4 and Figure 25 show the results using Table 1 condition except the latent dimension is increased to 5. For z_1 to z_3 , each value is close to Table 1. For z_4 and z_5 , $(2/\beta)\sigma_j^2 D'_j$ are almost 0 and the averages of $\sigma_{j(\mathbf{x})}^{-2}$ are close to 1. In such dimensions, $\text{Var}(y_j) < \beta/2$ and $D_{\text{KL}j}(\cdot) = 0$ will hold as explained in Appendix A.7, E.2, and B.6 (RD theory). Figure 25 describes the plot for $a_x^{n/2} \exp(-L_x/\beta)$ corresponding to Fig. 3d, also showing almost proportionality.

variable		z_1	z_2	z_3	z_4	z_5
$\frac{2}{\beta}\sigma_j^2 D'_j$	Av.	0.963	0.918	0.964	0.000	0.000
	SD	0.053	0.169	0.103	0.000	0.000
$\sigma_{j(\mathbf{x})}^{-2}$	Av.	3.34e1	1.46e2	5.88e2	1.00e0	1.00e0
	(Ratio) Av.	1.0	4.39	17.69	0.03	0.03

Table 4. Property measurements of the toy dataset with 5-dimensional latents trained using the square error loss.

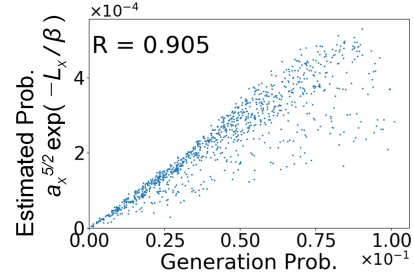


Figure 25. Plots of the data generation probability (x-axis) versus estimated probabilities (y-axes) for the square error loss. The dimension of latents is set to 5.

E. Additional results in CelebA dataset

E.1. Traversed outputs for all the component in the experimental section 5.2

Figure 26 shows decoder outputs for all the components, where each latent variable is traversed from -2 to 2 . The estimated variance of each y_j , i.e., σ_j^{-2} , is also shown in these figures. The latent variables z_i are numbered in descending order by the estimated variances. Figure 26a is a result using the conventional loss form, i.e., $L_{\mathbf{x}} = D(\mathbf{x}, \hat{\mathbf{x}}) + \beta D_{\text{KL}}(\cdot)$. The degrees of change seem to descend in accordance with the estimated variances. In the range where j is 1 from 10, the degrees of changes are large. In the range $j > 10$, the degrees of changes becomes gradually smaller. Furthermore, almost no change is observed in the range $j > 27$. As shown in Figure 4, $D_{\text{KL}(j)}(\cdot)$ is close to zero for $j > 27$, meaning no information. Note that the behavior of dimensional components where $D_{\text{KL}(j)}(\cdot) = 0$ is explained in section E.2. Thus, this result is clearly consistent with our theoretical analysis in section 4.3.

Figure 26b is a result using the decomposed loss form, i.e., $L_{\mathbf{x}} = D(\mathbf{x}, \check{\mathbf{x}}) + D(\check{\mathbf{x}}, \hat{\mathbf{x}}) + \beta D_{\text{KL}}(\cdot)$. The degrees of change also seem to descend in accordance with the estimated variances. When looking at the detail, there are still minor changes even $j = 32$. As shown in Figure 5, KL divergences $D_{\text{KL}(j)}(\cdot)$ for all the components are larger than zero. This implies all of the dimensional components have meaningful information. Therefore, we can see a minor change even $j = 32$. Thus, this result is also consistent with our theoretical analysis in Section 4.3.

Another minor difference is sharpness. Although the quantitative comparison is difficult, the decoded images in Figure 26b seems somewhat sharper than those in Figure 26a. A possible reason for this minor difference is as follows. The transform loss $D(\mathbf{x}, \check{\mathbf{x}})$ serves to bring the decoded image of $\mu_{(x)}$ closer to the input. In the conventional image coding, the orthonormal transform and its inverse transform are used for encoding and decoding, respectively. Therefore, the input and the decoded output are equivalent when not using quantization. If not so, the quality of the decoded image will suffer from the degradation. Considering this analogy, the use of decomposed loss might improve the decoded images for $\mu_{(x)}$, encouraging the improvement of the orthonormality of the encoder/decoder in VAE.

E.2. The understanding of latent components where $D_{\text{KL}(j)}(\cdot) = 0$ in Figure 4

This section explains the behaviors of latent components where $D_{\text{KL}(j)}(\cdot) = 0$, especially in Fig. 4. First, we explain why the norm becomes 0 when $D_{\text{KL}(j)}(\cdot) = 0$. The loss in Eq. 13 consists of a norm (multiplied by $\beta/2$) and $\beta D_{\text{KL}(j)}(\cdot)$ to find the best balance (trade-off) between them. If $D_{\text{KL}(j)}(\cdot) = 0$, the norm also becomes zero because balancing them is no more needed. Second, we explain the condition where $D_{\text{KL}(j)}(\cdot) = 0$. Let $\text{Var}(y_j)$ and $\sigma_{y_j(x)}^2$ be the variance and posterior variance of j -th implicit isometric component y_j , respectively. Here, $\sigma_{y_j(x)}^2$ is $\beta/2$ in our theory. Then the condition where $D_{\text{KL}(j)}(\cdot) = 0$ is derived as $\text{Var}(y_j) \leq \sigma_{y_j(x)}^2$, as shown in Appendix A.7 and B.6 (RD theory). In RD theory, this is corresponding to the case where the signal magnitude is always less than the quantizer size ($\sqrt{\beta/2}$ in β -VAE case) and no information is needed to be encoded. Finally, we explain the reason why the behaviors in Figs. 4 and 5 are different. In Fig. 5 with the decomposed loss, $\sigma_{y_j(x)}^2$ is almost $\beta/2$ as the theory expects. In this case, all of $\text{Var}(y_j)$ happen to be greater than $\sigma_{y_j(x)}^2$. In Fig. 4 with the conventional loss, however, $\sigma_{y_j(x)}^2$ is about 1.83 times greater than $\beta/2$. Note that in both Figs. 4 and 5, $\text{Var}(y_j)$ will be almost the same because of the isometric embedding. Since $\sigma_{y_j(x)}^2$ becomes larger, the number of dimensions where $\text{Var}(y_j) \leq \sigma_{y_j(x)}^2$ will increase. Accordingly, the dimensions where the norms are zero also increase. Figure 27 shows the CelebA results with smaller β , resulting smaller $\sigma_{y_j(x)}^2$. Here, all dimensions have nonzero norms because $\text{Var}(y_j) > \sigma_{y_j(x)}^2$ will hold.

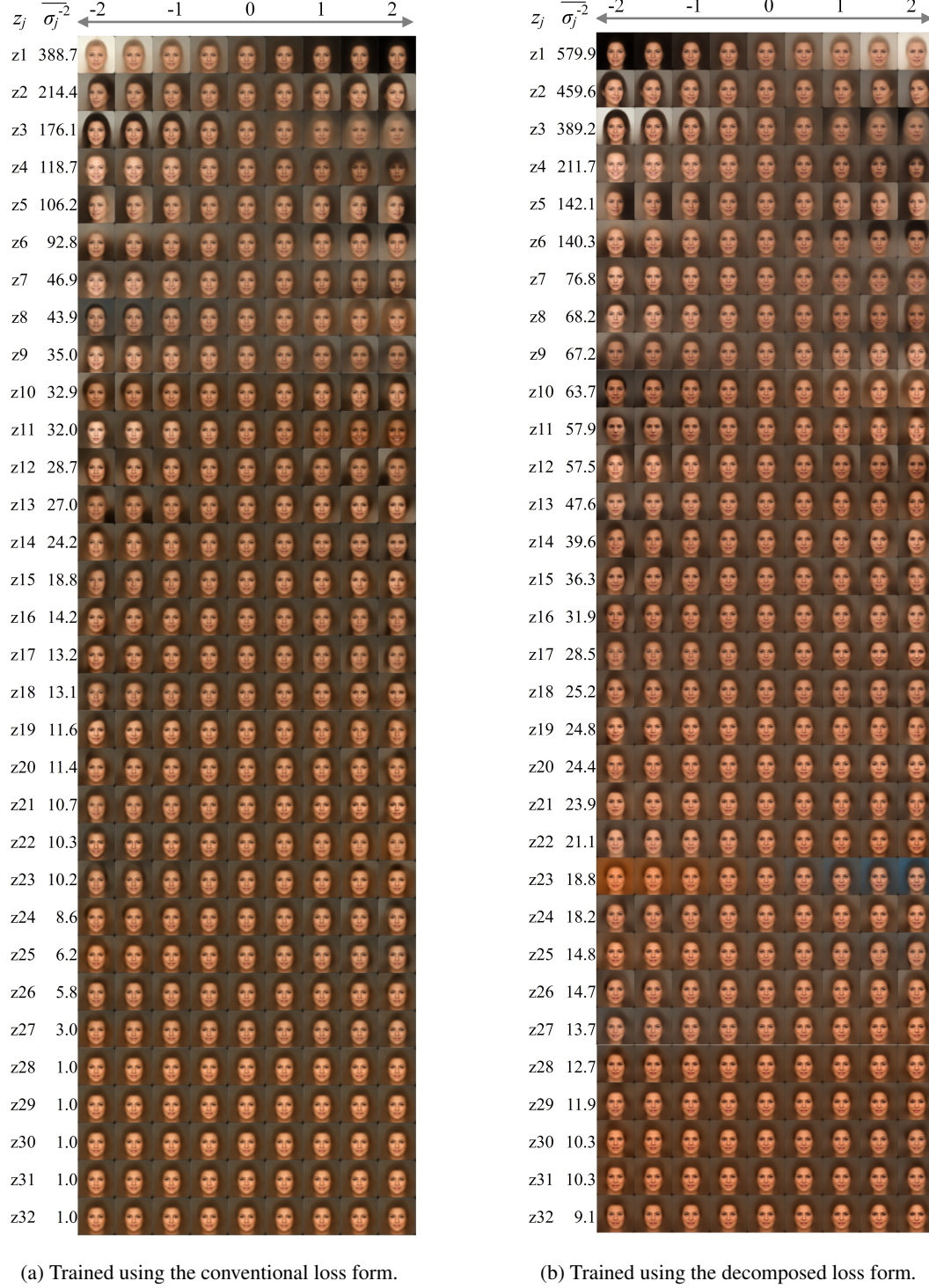


Figure 26. Traversed outputs for all the component, changing z_j from -2 to 2 . The latent variables z_j are numbered in descending order by the estimated variance σ_j^{-2} shown in Figures 4 and 5.

E.3. Additional experimental result with other condition

In this Section, we provide the experimental results with other condition. We use essentially the same condition as described in Appendix C.2, except for the following conditions. The bottleneck size and λ are set to 256 and 10000, respectively. The

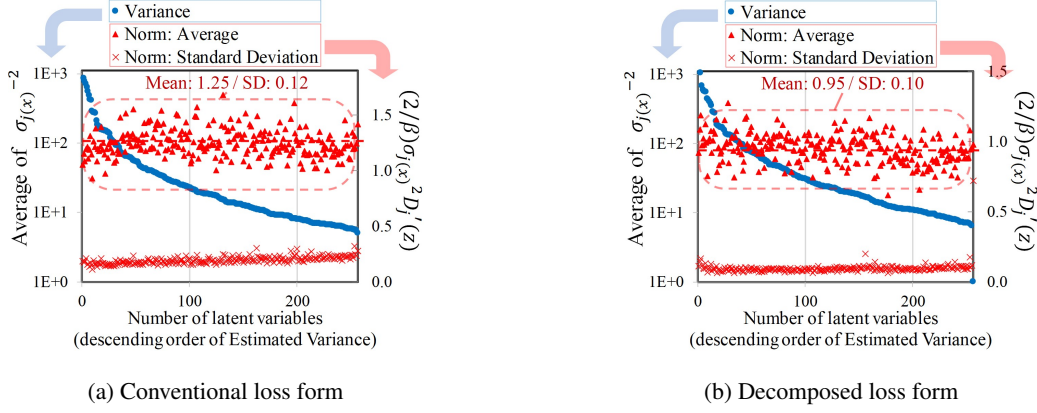


Figure 27. Graph of $\sigma_{j(x)}^{-2}$ average and $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$ in CelebA dataset. The bottleneck size and λ are set to 256 and 10000, respectively.

encoder network is composed of CNN(9, 9, 2, 64, GDN) - CNN(5, 5, 2, 64, GDN) - CNN(5, 5, 2, 64, GDN) - CNN(5, 5, 2, 64, GDN) - FC(1024, 2048, softplus) - FC(2048, 256, None) $\times 2$ (for μ and σ) in encoder. The decoder network is composed of FC(256, 2048, softplus) - FC(2048, 1024, softplus) - CNN(5, 5, 2, 64, IGDN) - CNN(5, 5, 2, 64, IGDN) - CNN(5, 5, 2, 64, IGDN)-CNN(9, 9, 2, 3, IGDN).

Figures 27a and 27b show the averages of $\sigma_{j(x)}^{-2}$ as well as the average and the standard deviation of $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$ in the conventional loss form and the decomposed loss form, respectively. When using the conventional loss form, the mean of $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$ is 1.25, which is closer to 1 than the mean 1.83 in Section 5.2. This suggests that the implicit transform is closer to the orthonormal. The possible reason is that a bigger reconstruction error is likely to cause the interference to RD-trade off and a slight violation of the theory, and it might be compensated with a larger lambda. When using the decomposed loss form, the mean of $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$ is 0.95, meaning almost unit norm. These results also support that VAE provides the implicit orthonormal transform even if the lambda or bottleneck size is varied.

F. Additional Experimental Result with MNIST dataset

In this Appendix, we provide the experimental result of Section 5.2 with MNIST dataset[‡] consists of binary hand-written digits with a dimension of 768(=28 \times 28). We use standard training split which includes 50,000 data points. For the reconstruction loss, we use the binary cross entropy loss (BCE) for the Bernoulli distribution. We averaged BCE by the number of pixels.

The encoder network is composed of FC(768, 1024, relu) - FC(1024, 1024, relu) - FC(1024, bottleneck size) in encoder. The decoder network is composed of FC(bottleneck size, 1024, relu) - FC(1024, 1024, relu) - FC(1024, 768, sigmoid). The batch size is 256 and the training iteration number is 50,000. In this section, results with two parameters, (bottleneck size=32, $\lambda=2000$) and (bottleneck size=64, $\lambda=10000$) are provided. Note that since we averaged BCE loss by the number of pixels, β in the conventional β VAE is derived by 768/ λ . Then, the model is optimized by Adam optimizer with the learning rate of 1e-3, using the conventional (not decomposed) loss form.

We use a PC with CPU Intel(R) Core(TM) i7-6850K CPU @ 3.60GHz, 12GB memory equipped with NVIDIA GeForce GTX 1080. The simulation time for each trial is about 10 minutes, including the statistics evaluation codes.

Figure 28 shows the averages of $\sigma_{j(x)}^{-2}$ as well as the average and the standard deviation of $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$. In both conditions, the means of $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$ averages are also close to 1 except in the dimensions where $\sigma_{j(x)}^{-2}$ is less than 10. These results suggest the theoretical property still holds when using the BCE loss. In the dimensions where $\sigma_{j(x)}^{-2}$ is less than 10, the $\frac{2}{\beta}\sigma_{j(x)}^2 D'_j(\mathbf{z})$ is somewhat lower than 1. The possible reason is that $D_{KL(j)}(\cdot)$ in such dimension is 0 for some inputs and is larger than 0 in other inputs. The understanding of the transition region needs further study.

[‡]<http://yann.lecun.com/exdb/mnist/>

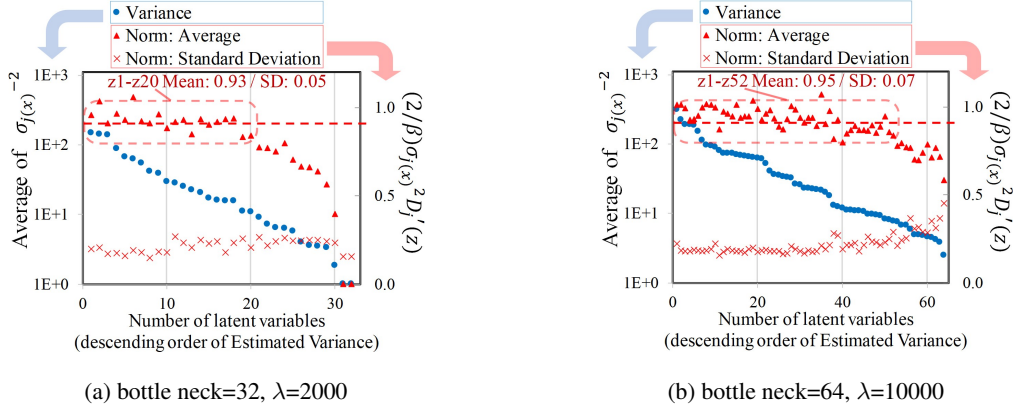


Figure 28. Graph of $\sigma_{j(\mathbf{x})}^{-2}$ average and $\frac{2}{\beta} \sigma_{j(\mathbf{x})}^2 D'_j(\mathbf{z})$ in MNIST dataset.

G. Derivation/Explanation in RDO-related equation expansions

G.1. Approximation of distortion in uniform quantization

Let T be a quantization step. Quantized values \hat{z}_j is derived as $k T$, where $k = \text{round}(z_j/T)$. Then d , the distortion per channel, is approximated by

$$\begin{aligned}
 d &= \sum_k \int_{(k-1/2)T}^{(k+1/2)T} p(z_j)(z_j - kT)^2 dz_j \simeq \sum_k T p(kT) \int_{(k-1/2)T}^{(k+1/2)T} \frac{1}{T} (z_j - kT)^2 dz_j \\
 &= \frac{T^2}{12} \sum_k T p(kT) \simeq \frac{T^2}{12}.
 \end{aligned} \tag{97}$$

Here, $\sum_k T p(kT) \simeq \int_{-\infty}^{\infty} p(z_j) dz_j = 1$ is used. The distortion for the given quantized value is also estimated as $T^2/12$, because this value is approximated by $\int_{(k-1/2)T}^{(k+1/2)T} \frac{1}{T} (z_j - kT)^2 dz_j$.

G.2. Approximation of reconstruction loss as a quadratic form.

In this appendix, the approximations of the reconstruction losses as a quadratic form ${}^t \delta \mathbf{x} \mathbf{G}_{\mathbf{x}} \delta \mathbf{x} + C_{\mathbf{x}}$ are explained for the sum of square error (SSE), binary cross entropy (BCE) and Structural Similarity (SSIM). Here, we have borrowed the derivation of BCE and SSIM from Kato et al. (2020), and add some explanation and clarification to them for convenience. We also describe the log-likelihood of the Gaussian distribution.

Let $\hat{\mathbf{x}}$ and \hat{x}_i be decoded sample $\text{Dec}_{\theta}(\mathbf{z})$ and its i -th dimensional component respectively. $\delta \mathbf{x}$ and δx_i denote $\mathbf{x} - \hat{\mathbf{x}}$ and $x_i - \hat{x}_i$, respectively. It is also assumed that $\delta \mathbf{x}$ and δx_i are infinitesimal. The details of the approximations are described as follows.

Sum square error:

In the case of sum square error, $\mathbf{G}_{\mathbf{x}}$ is equal to \mathbf{I}_m . This can be derived as:

$$\sum_{i=1}^m (x_i - \hat{x}_i)^2 = \sum_{i=1}^m \delta x_i^2 = {}^t \delta \mathbf{x} \mathbf{I}_m \delta \mathbf{x}. \tag{98}$$

Binary cross entropy:

Binary cross entropy is a log likelihood of the Bernoulli distribution. The Bernoulli distribution is described as:

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^m \hat{x}_i^{x_i} (1 - \hat{x}_i)^{(1-x_i)}. \tag{99}$$

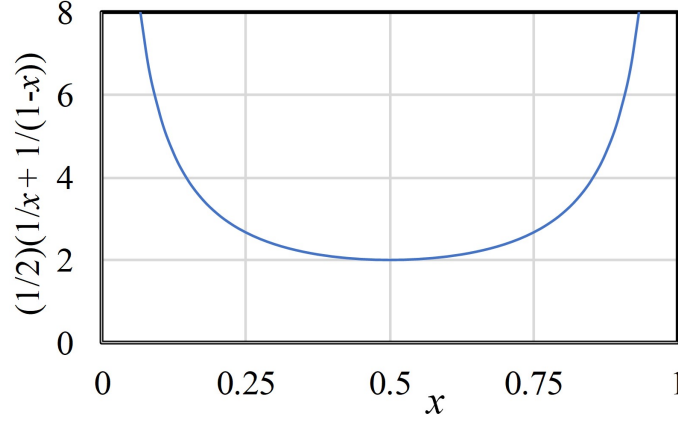


Figure 29. Graph of $\frac{1}{2} \left(\frac{1}{x} + \frac{1}{1-x} \right)$ in the BCE approximation.

Then, the binary cross-entropy (BCE) can be expanded as:

$$\begin{aligned}
 -\log p_{\theta}(\mathbf{x}|\mathbf{z}) &= -\log \prod_{i=1}^m \hat{x}_i^{x_i} (1 - \hat{x}_i)^{(1-x_i)} \\
 &= \sum_{i=1}^m (-x_i \log \hat{x}_i - (1 - x_i) \log (1 - \hat{x}_i)) \\
 &= \sum_i \left(-x_i \log \left(1 + \frac{\delta x_i}{x_i} \right) - (1 - x_i) \log \left(1 - \frac{\delta x_i}{1 - x_i} \right) \right) \\
 &\quad + \sum_i (-x_i \log(x_i) - (1 - x_i) \log(1 - x_i)).
 \end{aligned} \tag{100}$$

Here, the second term of the last equation is a constant $C_{\mathbf{x}}$ depending on \mathbf{x} . Using $\log(1 + x) = x - x^2/2 + O(x^3)$, the first term of the last equation is further expanded as follows:

$$\begin{aligned}
 \sum_i \left(-x_i \left(\frac{\delta x_i}{x_i} - \frac{\delta x_i^2}{2x_i^2} \right) - (1 - x_i) \left(-\frac{\delta x_i}{1 - x_i} - \frac{\delta x_i^2}{2(1 - x_i)^2} \right) + O(\delta x_i^3) \right) \\
 = \sum_i \left(\frac{1}{2} \left(\frac{1}{x_i} + \frac{1}{1 - x_i} \right) \delta x_i^2 + O(\delta x_i^3) \right).
 \end{aligned} \tag{101}$$

As a result, a metric tensor $\mathbf{G}_{\mathbf{x}}$ can be approximated as the following positive definite Hermitian matrix:

$$\mathbf{G}_{\mathbf{x}} = \begin{pmatrix} \frac{1}{2} \left(\frac{1}{x_1} + \frac{1}{1-x_1} \right) & 0 & \dots \\ 0 & \frac{1}{2} \left(\frac{1}{x_2} + \frac{1}{1-x_2} \right) & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}. \tag{102}$$

Here, the loss function in each dimension $\frac{1}{2} \left(\frac{1}{x_1} + \frac{1}{1-x_1} \right)$ is a downward-convex function as shown in Figure 29.

Structural similarity (SSIM):

Structural similarity (SSIM) (Wang et al., 2001) is widely used for picture quality metric, which is close to subjective quality. Let SSIM be a SSIM value between two pictures. The range of the SSIM is between 0 and 1. The higher the value, the better the quality. In this appendix, we also show that $(1 - \text{SSIM})$ can be approximated to a quadratic form such as ${}^t \delta \mathbf{x} \mathbf{G}_{\mathbf{x}} \delta \mathbf{x}$.

$\text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{y})$ denotes a SSIM value between $N \times N$ windows in pictures X and Y , where $\mathbf{x} \in \mathbb{R}^{N^2}$ and $\mathbf{y} \in \mathbb{R}^{N^2}$ denote $N \times N$ pixels cropped from the top-left coordinate (h, v) in the images X and Y , respectively. Let $\mu_{\mathbf{x}}, \mu_{\mathbf{y}}$ be the

averages of all dimensional components in \mathbf{x} , \mathbf{y} , and $\sigma_{\mathbf{x}}$, $\sigma_{\mathbf{y}}$ be the variances of all dimensional components in \mathbf{x} , \mathbf{y} in the $N \times N$ windows, respectively. Then, $\text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{y})$ is derived as

$$\text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{y}) = \frac{2\mu_{\mathbf{x}}\mu_{\mathbf{y}}}{\mu_{\mathbf{x}}^2 + \mu_{\mathbf{y}}^2} \cdot \frac{2\sigma_{\mathbf{x}\mathbf{y}}}{\sigma_{\mathbf{x}}^2 + \sigma_{\mathbf{y}}^2}. \quad (103)$$

In order to calculate a SSIM value for a picture, the window is shifted in a whole picture and all of SSIM values are averaged. Therefore, if $(1 - \text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{y}))$ is expressed as a quadratic form ${}^t\delta\mathbf{x} \mathbf{G}_{(h,v)\mathbf{x}} \delta\mathbf{x}$, $(1 - \text{SSIM})$ can be also expressed in quadratic form ${}^t\delta\mathbf{x} \mathbf{G}_{\mathbf{x}} \delta\mathbf{x}$.

Let $\delta\mathbf{x}$ be a minute displacement of \mathbf{x} . $\mu_{\delta\mathbf{x}}$ and $\sigma_{\delta\mathbf{x}}^2$ denote an average and variance of all dimensional components in $\delta\mathbf{x}$, respectively. Then, SSIM between \mathbf{x} and $\mathbf{x} + \delta\mathbf{x}$ can be approximated as:

$$\text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{x} + \delta\mathbf{x}) \simeq 1 - \frac{\mu_{\delta\mathbf{x}}^2}{2\mu_{\mathbf{x}}^2} - \frac{\sigma_{\delta\mathbf{x}}^2}{2\sigma_{\mathbf{x}}^2} + O\left((|\delta\mathbf{x}|/|\mathbf{x}|)^3\right). \quad (104)$$

Then $\mu_{\delta\mathbf{x}}^2$ and $\sigma_{\delta\mathbf{x}}^2$ can be expressed as

$$\mu_{\delta\mathbf{x}}^2 = {}^t\delta\mathbf{x} \mathbf{M} \delta\mathbf{x}, \quad \text{where } \mathbf{M} = \frac{1}{N^2} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}, \quad (105)$$

and

$$\sigma_{\delta\mathbf{x}}^2 = {}^t\delta\mathbf{x} \mathbf{V} \delta\mathbf{x}, \quad \text{where } \mathbf{V} = \frac{1}{N} \mathbf{I}_N - \mathbf{M}, \quad (106)$$

respectively. As a result, $(1 - \text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{x} + \delta\mathbf{x}))$ can be expressed in the following quadratic form as:

$$1 - \text{SSIM}_{N \times N(h,v)}(\mathbf{x}, \mathbf{x} + \delta\mathbf{x}) \simeq {}^t\delta\mathbf{x} \mathbf{G}_{(h,v)\mathbf{x}} \delta\mathbf{x}, \quad \text{where } \mathbf{G}_{(h,v)\mathbf{x}} = \left(\frac{1}{2\mu_{\mathbf{x}}^2} \mathbf{M} + \frac{1}{2\sigma_{\mathbf{x}}^2} \mathbf{V} \right). \quad (107)$$

It is noted that \mathbf{M} is a positive definite Hermitian matrix and \mathbf{V} is a positive semidefinite Hermitian matrix. Therefore, $\mathbf{G}_{(h,v)\mathbf{x}}$ is a positive definite Hermitian matrix. As a result, $(1 - \text{SSIM})$ can be also expressed in quadratic form ${}^t\delta\mathbf{x} \mathbf{G}_{\mathbf{x}} \delta\mathbf{x}$, where $\mathbf{G}_{\mathbf{x}}$ is a positive definite Hermitian matrix.

Log-likelihood of Gaussian distribution:

Gaussian distribution is described as:

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - \hat{x}_i)^2/2\sigma^2} = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\delta x_i^2/2\sigma^2}, \quad (108)$$

where σ^2 is a variance as a hyper parameter. Then, the log-likelihood of the Gaussian distribution is denoted as:

$$-\log p_{\theta}(\mathbf{x}|\mathbf{z}) = -\log \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\delta x_i^2/2\sigma^2} = \frac{1}{2\sigma^2} \sum_{i=1}^m \delta x_i^2 + \frac{m}{2} \log(2\pi\sigma^2). \quad (109)$$

Since the first term is $(1/2\sigma^2) {}^t\delta\mathbf{x} \mathbf{I}_m \delta\mathbf{x}$, $\mathbf{G}_{\mathbf{x}} = (1/2\sigma^2) \mathbf{I}_m$ holds. $C_{\mathbf{x}}$ is the second term of the last equation in Eq.109.

H. Detail of the Experiment in Section 5.3

In this section, we provide further detail of experiment in Section 5.3.

H.1. Datasets

We describe the detail of following four public datasets:

Table 5. Hyper parameter for RaDOGAGA

Dataset	Autoencoder	Transform loss	λ_1	λ_2	
KDDCup99	200, 100, 10, 100, 200	L2	30000	6000	
Thyroid	60, 30, 6, 30, 60	L2	6000	18000	
Arrhythmia	200, 100, 50, 100, 200	L2	6000	24000	
KDDCup-rev	200, 50, 20, 50, 200	SSE	30000	6000	

KDDCUP99 (Dua & Graff, 2019) The KDDCUP99 10 percent dataset from the UCI repository is a dataset for cyber-attack detection. This dataset includes 494,021 instances. Each instance contains 34 continuous and 7 categorical features. We use one hot representation to encode the categorical features, and finally obtain a dataset with features of 121 dimensions. Only 20% of instances labeled -normal- and the rest labeled as -attacks-. Therefore, -normal- instances are used as anomalies, because they are in a minority group.

Thyroid (Dua & Graff, 2019) This dataset consists of 3,772 data sample with 6-dimensional feature from patients. Each instance can be divided in three classes: normal (not hypothyroid), hyperfunction, and subnormal functioning. We regard the hyperfunction class (2.5%) as an anomaly and rest two classes as normal.

Arrhythmia (Dua & Graff, 2019) This is dataset to detect cardiac arrhythmia which containing 452 instances with 274-dimensional feature. We treat minor classes (3, 4, 5, 7, 8, 9, 14, and 15, accounting for 15% of the total) as anomalies. The rest of classes are treated as normal.

KDDCUP-Rev (Dua & Graff, 2019) This is a revised version of KDDCUP99. To treat “normal” instances as the majority in the KDDCUP dataset, we keep all -normal- instances and randomly pick up -attack- instances so that the ratio of -normal- and -attack- to be 8:2. The number of instances is 121,597 in the end.

Data is max-min normalized toward dimension through the entire dataset, which is the same setting as previous studies.

H.2. Network architecture, Hyperparameter, and Training Detail

The VAE in this experiment consists of FC layers. Expect for the last layer of the encoder, Leaky ReLU (for KDDCup99, Thyroid, and Arrhythmia) or tanh (for KDDCup-rev) is attached as the activation function.

In this experiment, VAE is constructed by the form of decomposed loss to promote isometricity as explained in Remark 1. Here, the decomposed loss function L_x is set to $\lambda_1 D(x, \check{x}) + \lambda_2 D(\check{x}, \hat{x}) + D_{KL}(\cdot)$, meaning, $\lambda_2 = \beta^{-1}$, are adjusted independently for reconstruction loss and transform loss. For the transform loss, we tested both L2 norm and SSE loss and choose the better one for each dataset. The reason for introducing L2 loss to the transform loss is as follows. The reduction of the transform loss promotes the isometricity, as explained in Remark 1 of Section 4.2. Since the derivative of L2 norm is steeper than SSE used in coding loss, the use of L2 norm for transform loss will reduce the value of transform loss explicitly and promote the isometricity.

Hyperparameter is described in Table 5. The first column is the number of neurons. For Thyroid, we also tested the size of (30, 24, 6, 24, 30). For other datasets, we tested the size of (200, 100 or 50, 10 or 20 or 50, 100 or 50, 200). The second column is the type of reconstruction loss. (λ_1, λ_2) is determined experimentally. Both of them varied from 6000 to 30000 by every 6000 intervals. For all datasets, optimization is done by Adam optimizer with a learning rate of 1×10^{-4} with batch size of 1024. The epoch numbers for each dataset are 600, 40000, 30000, and 600 respectively. Test models are saved by every 1/10 epochs and early stop is applied. For this experiment, we use GeForce GTX 1080.

H.3. Precision, Recall, and F1

Due to the page limitation, we reported only F1 score in main paper. Now we provide Precision and Recall Score as well in Table 6.

Table 6. Average and standard deviations (in brackets) of Precision, Recall and F1

Dataset	Methods	Precision	Recall	F1
KDDCup	GMVAE [†]	0.952	0.9141	0.9326
	DAGMM [†]	0.9427 (0.0052)	0.9575 (0.0053)	0.9500 (0.0052)
	RaDOGAGA(d) [†]	0.9550 (0.0037)	0.9700 (0.0038)	0.9624 (0.0038)
	RaDOGAGA(log(d)) [†]	0.9563 (0.0042)	0.9714 (0.0042)	0.9638 (0.0042)
	VAE	0.9568(0.0007)	0.9718 (0.0007)	0.9642(0.0007)
Thyroid	GMVAE [†]	0.7105	0.5745	0.6353
	DAGMM [†]	0.4656 (0.0481)	0.4859 (0.0502)	0.4755 (0.0491)
	RaDOGAGA(d) [†]	0.6313 (0.0476)	0.6587 (0.0496)	0.6447 (0.0486)
	RaDOGAGA(log(d)) [†]	0.6562 (0.0572)	0.6848 (0.0597)	0.6702 (0.0585)
	VAE	0.6458 (0.04270)	0.6739 (0.04455)	0.6596 (0.0436)
Arrhythmia	GMVAE [†]	0.4375	0.4242	0.4308
	DAGMM [†]	0.4985 (0.0389)	0.5136 (0.0401)	0.5060 (0.0395)
	RaDOGAGA(d) [†]	0.5353 (0.0461)	0.5515 (0.0475)	0.5433 (0.0468)
	RaDOGAGA(log(d)) [†]	0.5294 (0.0405)	0.5455 (0.0418)	0.5373 (0.0411)
	VAE	0.4912(0.0406)	0.5061 (0.0419)	0.4985 (0.0413)
KDDCup-rev	DAGMM [†]	0.9778 (0.0018)	0.9779 (0.0017)	0.9779 (0.0018)
	RaDOGAGA(d) [†]	0.9768 (0.0033)	0.9827 (0.0012)	0.9797 (0.0015)
	RaDOGAGA(log(d)) [†]	0.9864 (0.0009)	0.9865 (0.0009)	0.9865 (0.0009)
	VAE	0.9880 (0.0008)	0.9881 (0.0008)	0.9880 (0.0008)

[†]Scores are cited from Liao et al. (2018) (GMVAE) and Kato et al. (2020)(DAGMM, RaDOGAGA)