

A. How Jacobian Matrix Forms a Constantly Scaled Orthonormal System

In this appendix, we derive equations corresponding to Eqs. (13) and (14) for the case of $M > N$. The guiding principle of derivation is the same as in Section 4.2: examining the condition to minimize the expected loss. As in Section 4.2, we assume that the encoder and the decoder are trained enough in terms of reconstruction error so that $\mathbf{x} \simeq \hat{\mathbf{x}}$ holds and the second term $\lambda_1 h(D(\mathbf{x}, \hat{\mathbf{x}}))$ in Eq. (7) can be ignored.

We assume that the Jacobian matrix $\mathbf{J}(\mathbf{z}) = \partial \mathbf{x} / \partial \mathbf{z} = \partial g_\phi(\mathbf{z}) / \partial \mathbf{z} \in \mathbb{R}^{M \times N}$ is full-rank at every point $\mathbf{z} \in \mathbb{R}^N$ as in Section 4.2. Based on Eq. (6), Eq. (4) and Taylor expansion, the difference $\check{\mathbf{x}} - \hat{\mathbf{x}}$ can be approximated by $\epsilon = \sum_{i=1}^N \epsilon_i (\partial \mathbf{x} / \partial z_i) \in \mathbb{R}^M$. As in Section 4.2, the expectation of the third term in Eq. (7) is re-written as follows:

$$E_{\epsilon \sim P_\epsilon(\epsilon)} [\epsilon^\top \mathbf{A}(\mathbf{x}) \epsilon] = \sigma^2 \sum_{j=1}^N \left(\frac{\partial \mathbf{x}}{\partial z_j} \right)^\top \mathbf{A}(\mathbf{x}) \left(\frac{\partial \mathbf{x}}{\partial z_j} \right). \quad (20)$$

This equation has the same form as Eq. (10) except the differences in dimensions: $\partial \mathbf{x} / \partial z_j \in \mathbb{R}^M$ and $\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{M \times M}$ in Eq. (20) while $\partial \mathbf{x} / \partial z_j \in \mathbb{R}^N$ and $\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{N \times N}$ in Eq. (10). We have essentially no difference from Section 4.2 so far.

However, from this point, we cannot follow the same way we used in Section 4.2 to derive the equation corresponding to Eq. (13), due to the mismatch of M and N . Yet, as we show below, step-by-step modifications lead us to the same conclusion.

Firstly, note that we can always regard g_ϕ as a composition function by inserting a smooth invertible function $\rho : \mathbb{R}^N \rightarrow \mathbb{R}^N$ and its inverse as follows:

$$g_\phi(\mathbf{z}) = g_\phi(\rho^{-1}(\rho(\mathbf{z}))) = \tilde{g}_\phi(\rho(\mathbf{z})). \quad (21)$$

Let $\mathbf{y} \in \mathbb{R}^N$ be an auxiliary variable defined by $\mathbf{y} = \rho(\mathbf{z})$. Due to the chain rule of differentiation, $\partial \mathbf{x} / \partial \mathbf{z}$ can be represented as

$$\frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \mathbf{G} \mathbf{B}, \quad (22)$$

where we define \mathbf{G} and \mathbf{B} as $\mathbf{G} = \partial \mathbf{x} / \partial \mathbf{y} \in \mathbb{R}^{M \times N}$, $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_N) = \partial \mathbf{z} / \partial \mathbf{y} \in \mathbb{R}^{N \times N}$.

Since \mathbf{z} and \mathbf{y} have the same dimension N , the relationship between $P_z(\mathbf{z})$ and $P_y(\mathbf{y})$ is described by $|\det(\mathbf{B})|$ (the absolute value of Jacobian determinant), which corresponds to the volume change under the function ρ , as follows:

$$P_z(\mathbf{z}) = |\det(\mathbf{B})| P_y(\mathbf{y}). \quad (23)$$

Thus the expectation of L in Eq. (7) can be approximated as follows:

$$E_{\epsilon \sim P_\epsilon(\epsilon)} [L] \simeq -\log(|\det(\mathbf{B})|) - \log(P_y(\mathbf{y})) + \lambda_2 \sigma^2 \left(\sum_{j=1}^N (\mathbf{G} \mathbf{b}_j)^\top \mathbf{A}(\mathbf{x}) (\mathbf{G} \mathbf{b}_j) \right). \quad (24)$$

To derive the minimization condition of the expected loss, we need further preparations. Let \tilde{b}_{ij} denote the cofactor of matrix \mathbf{B} with regard to the element b_{ij} . We define a vector $\tilde{\mathbf{b}}_j$ as follows:

$$\tilde{\mathbf{b}}_j = \begin{pmatrix} \tilde{b}_{1j} \\ \tilde{b}_{2j} \\ \vdots \\ \tilde{b}_{Nj} \end{pmatrix}. \quad (25)$$

The following equation is a property of the cofactor (Strang, 2006):

$$\mathbf{b}_i^\top \tilde{\mathbf{b}}_j = \sum_{k=1}^N b_{ki} \tilde{b}_{kj} = \delta_{ij} \det(\mathbf{B}). \quad (26)$$

In addition, since $|\det(\mathbf{B})| = (\det(\mathbf{B})^2)^{\frac{1}{2}} = ((\sum_{k=1}^N b_{kj}\tilde{b}_{kj})^2)^{\frac{1}{2}}$, we have the following result:

$$\frac{\partial |\det(\mathbf{B})|}{\partial b_{ij}} = \frac{1}{2} \left(\sum_{k=1}^N b_{kj}\tilde{b}_{kj} \right)^{-\frac{1}{2}} \cdot 2 \left(\sum_{k=1}^N b_{kj}\tilde{b}_{kj} \right) \tilde{b}_{ij} = \frac{\det(\mathbf{B})}{|\det(\mathbf{B})|} \tilde{b}_{ij}. \quad (27)$$

Therefore, the following equations hold:

$$\frac{\partial \log(|\det(\mathbf{B})|)}{\partial b_{ij}} = \frac{1}{|\det(\mathbf{B})|} \frac{\partial |\det(\mathbf{B})|}{\partial b_{ij}} = \frac{1}{|\det(\mathbf{B})|} \frac{\det(\mathbf{B})}{|\det(\mathbf{B})|} \tilde{b}_{ij} = \frac{\det(\mathbf{B})}{\det(\mathbf{B})^2} \tilde{b}_{ij} = \frac{1}{\det(\mathbf{B})} \tilde{b}_{ij}, \quad (28)$$

$$\frac{\partial \log(|\det(\mathbf{B})|)}{\partial \mathbf{b}_j} = \frac{1}{\det(\mathbf{B})} \tilde{\mathbf{b}}_j. \quad (29)$$

By differentiating the right hand side of Eq. (24) by \mathbf{b}_j and setting the result to zero, the following equation is derived as a condition to minimize the expected loss:

$$2\lambda_2\sigma^2 \mathbf{G}^\top \mathbf{A}(\mathbf{x}) \mathbf{G} \mathbf{b}_j = \frac{1}{\det(\mathbf{B})} \tilde{\mathbf{b}}_j. \quad (30)$$

Here we used Eq. (29). By multiplying \mathbf{b}_i^\top to this equation from the left and dividing the result by $2\lambda_2\sigma^2$, we have

$$(\mathbf{G} \mathbf{b}_i)^\top \mathbf{A}(\mathbf{x}) (\mathbf{G} \mathbf{b}_j) = \frac{1}{2\lambda_2\sigma^2} \frac{1}{\det(\mathbf{B})} \mathbf{b}_i^\top \tilde{\mathbf{b}}_j \quad (31)$$

$$= \frac{1}{2\lambda_2\sigma^2} \delta_{ij}, \quad (32)$$

where the second line follows from Eq. (26). Since $\mathbf{G} \mathbf{b}_i = (\partial \mathbf{x} / \partial \mathbf{y})(\partial \mathbf{y} / \partial z_i) = \partial \mathbf{x} / \partial z_i$ and $\mathbf{G} \mathbf{b}_j = (\partial \mathbf{x} / \partial \mathbf{y})(\partial \mathbf{y} / \partial z_j) = \partial \mathbf{x} / \partial z_j$, we can come back to the expression with the original variables \mathbf{x} and \mathbf{z} and reach the following conclusion:

$$\left(\frac{\partial \mathbf{x}}{\partial z_i} \right)^\top \mathbf{A}(\mathbf{x}) \left(\frac{\partial \mathbf{x}}{\partial z_j} \right) = \frac{1}{2\lambda_2\sigma^2} \delta_{ij}. \quad (33)$$

Here the dimensions are different from Eq. (13) ($\partial \mathbf{x} / \partial z_i, \partial \mathbf{x} / \partial z_j \in \mathbb{R}^M$ and $\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{M \times M}$) but the meaning is same: the columns of the Jacobian matrix of two spaces $\partial \mathbf{x} / \partial z_1, \dots, \partial \mathbf{x} / \partial z_N$ form a constantly-scaled orthonormal system with respect to the inner product defined by $\mathbf{A}(\mathbf{x})$ at every point.

Now we can derive the second conclusion in the exactly same manner as in Section 4.2, although the dimensions are different ($\mathbf{v}_x, \mathbf{w}_x \in \mathbb{R}^M$, $\mathbf{A}(\mathbf{x}) \in \mathbb{R}^{M \times M}$ and $\mathbf{v}_z, \mathbf{w}_z \in \mathbb{R}^N$):

$$\begin{aligned} \mathbf{v}_x \mathbf{A}(\mathbf{x}) \mathbf{w}_x &= \sum_{i=0}^N \sum_{j=0}^N \left(\frac{\partial \mathbf{x}}{\partial z_i} v_{zi} \right)^\top \mathbf{A}(\mathbf{x}) \left(\frac{\partial \mathbf{x}}{\partial z_j} w_{zj} \right) \\ &= \frac{1}{2\lambda_2\sigma^2} \sum_{i=0}^N v_{zi} w_{zi} = \frac{1}{2\lambda_2\sigma^2} \mathbf{v}_z \cdot \mathbf{w}_z, \end{aligned} \quad (34)$$

which means the map is isometric in the sense of Eq. (2).

B. Product of Singular Values as a Generalization of the Absolute Value of Jacobian Determinant

In this appendix, we show the following two arguments we stated in the last part of Section 4.2: i) when a region in \mathbb{R}^N is mapped to \mathbb{R}^M by the decoder function, the ratio of the volume of the original region and its corresponding value is equal to the product of singular values of Jacobian matrix, ii) this quantity can be expressed by $\mathbf{A}(\mathbf{x})$ under a certain condition. The Jacobian matrix $\mathbf{J}(\mathbf{z}) = \partial \mathbf{x} / \partial \mathbf{z} = \partial g_\phi(\mathbf{z}) / \partial \mathbf{z} \in \mathbb{R}^{M \times N}$ is assumed to be full-rank as in Section 4.2 and Appendix A.

Let's consider the singular value decomposition $\mathbf{J}(\mathbf{z}) = \mathbf{U}(\mathbf{z}) \mathbf{\Sigma}(\mathbf{z}) \mathbf{V}(\mathbf{z})^\top$, where $\mathbf{U}(\mathbf{z}) \in \mathbb{R}^{M \times M}$, $\mathbf{\Sigma}(\mathbf{z}) \in \mathbb{R}^{M \times N}$, $\mathbf{V}(\mathbf{z}) \in \mathbb{R}^{N \times N}$. Note that $\{\mathbf{V}_{:,j}(\mathbf{z})\}_{j=1}^N$ is an orthonormal basis of \mathbb{R}^N and $\{\mathbf{U}_{:,j}(\mathbf{z})\}_{j=1}^M$ is an orthonormal basis of \mathbb{R}^M with respect to the standard inner product.

Consider an N -dimensional hypercube \mathbf{c} specified by $\{\varepsilon \mathbf{V}_{:,j}(\mathbf{z})\}_{j=1}^N$ ($\varepsilon > 0$) attached to $\mathbf{z} \in \mathbb{R}^N$. When ε is small, the effect of the decoder function on $\{\varepsilon \mathbf{V}_{:,j}(\mathbf{z})\}_{j=1}^N$ is approximated by $\mathbf{J}(\mathbf{z}) = \mathbf{U}(\mathbf{z})\boldsymbol{\Sigma}(\mathbf{z})\mathbf{V}(\mathbf{z})^\top$ and thus the mapped region of \mathbf{c} in \mathbb{R}^M is approximated by a region $\tilde{\mathbf{c}}$ specified by $\{\varepsilon \mathbf{J}(\mathbf{z})\mathbf{V}_{:,j}(\mathbf{z})\}_{j=1}^N = \{\varepsilon s_j(\mathbf{z})\mathbf{U}_{:,j}(\mathbf{z})\}_{j=1}^N$, where $s_1(\mathbf{z}) \geq \dots \geq s_N(\mathbf{z}) > 0$ are the singular values of $\mathbf{J}(\mathbf{z})$ (remember full-rank assumption we posed).

Therefore, while the volume of the original hypercube \mathbf{c} is ε^N , the corresponding value of $\tilde{\mathbf{c}} \in \mathbb{R}^M$ is $\varepsilon^N J_{sv}(\mathbf{z})$, where we define $J_{sv}(\mathbf{z})$ as $J_{sv}(\mathbf{z}) = s_1(\mathbf{z}) \cdots s_N(\mathbf{z})$, that is, the product of the singular values of the Jacobian matrix $\mathbf{J}(\mathbf{z})$. This relationship holds for any $\mathbf{z} \in \mathbb{R}^N$ and we can take arbitrary small ε . Thus, the ratio of the volume of an arbitrary region in \mathbb{R}^N and its corresponding value in \mathbb{R}^M is also $J_{sv}(\mathbf{z})^*$.

Let's move to the second argument. Note that Eq. (33) can be rewritten in the following form since $\mathbf{J}(\mathbf{z}) = (\partial \mathbf{x} / \partial z_1, \dots, \partial \mathbf{x} / \partial z_N)$:

$$\mathbf{J}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{J}(\mathbf{z}) = \frac{1}{2\lambda_2\sigma^2} \mathbf{I}_N. \quad (35)$$

Let $0 < \alpha_1(\mathbf{A}(\mathbf{x})) \leq \dots \leq \alpha_N(\mathbf{A}(\mathbf{x})) \leq \dots \leq \alpha_M(\mathbf{A}(\mathbf{x}))$ be the eigenvalues of $\mathbf{A}(\mathbf{x})$. If the condition

$$[\mathbf{O}_{(M-N) \times N} \quad \mathbf{I}_N] \mathbf{U}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{U}(\mathbf{z}) \begin{bmatrix} \mathbf{I}_N \\ \mathbf{O}_{(M-N) \times N} \end{bmatrix} = \mathbf{O}_{(M-N) \times N} \quad (36)$$

holds for all $\mathbf{z} \in \mathbb{R}^N$, the following relation holds for $J_{sv}(\mathbf{z})$:

$$J_{sv}(\mathbf{z}) = \left(\frac{1}{2\lambda_2\sigma^2} \right)^{\frac{N}{2}} \left(\alpha_1(\mathbf{A}(\mathbf{x})) \cdots \alpha_N(\mathbf{A}(\mathbf{x})) \right)^{-\frac{1}{2}}. \quad (37)$$

Here $\mathbf{O}_{(M-N) \times N} \in \mathbb{R}^{(M-N) \times N}$ denotes the matrix consisting of zeros.

To see this, let us first define $\mathbf{S}(\mathbf{z}) \in \mathbb{R}^{N \times N}$ as $\mathbf{S}(\mathbf{z}) = \text{diag}(s_1(\mathbf{z}), \dots, s_N(\mathbf{z}))$. Then $\mathbf{J}(\mathbf{z}) = \mathbf{U}(\mathbf{z})^\top [\mathbf{S}(\mathbf{z}) \mathbf{O}_{(M-N) \times N}]^\top \mathbf{V}(\mathbf{z})$. We obtain the following equation by substituting this expression of $\mathbf{J}(\mathbf{z})$ to Eq. (35):

$$\mathbf{V}(\mathbf{z}) [\mathbf{S}(\mathbf{z}) \mathbf{O}_{(M-N) \times N}] \mathbf{U}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{U}(\mathbf{z}) \begin{bmatrix} \mathbf{S}(\mathbf{z}) \\ \mathbf{O}_{(M-N) \times N} \end{bmatrix} \mathbf{V}(\mathbf{z})^\top = \frac{1}{2\lambda_2\sigma^2} \mathbf{I}_N. \quad (38)$$

Furthermore, we get the following equation by multiplying Eq. (38) by $\mathbf{S}(\mathbf{z})^{-1} \mathbf{V}(\mathbf{z})^\top$ from the left and $\mathbf{V}(\mathbf{z}) \mathbf{S}(\mathbf{z})^{-1}$ from the right:

$$[\mathbf{I}_N \quad \mathbf{O}_{(M-N) \times N}] \mathbf{U}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{U}(\mathbf{z}) \begin{bmatrix} \mathbf{I}_N \\ \mathbf{O}_{(M-N) \times N} \end{bmatrix} = \frac{1}{2\lambda_2\sigma^2} \mathbf{S}(\mathbf{z})^{-2} \quad (39)$$

This means $\mathbf{U}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{U}(\mathbf{z})$ has the following form:

$$\mathbf{U}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{U}(\mathbf{z}) = \begin{bmatrix} \frac{1}{2\lambda_2\sigma^2} \mathbf{S}(\mathbf{z})^{-2} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{D} \end{bmatrix}, \quad (40)$$

where $\mathbf{C} \in \mathbb{R}^{N \times (M-N)}$ and $\mathbf{D} \in \mathbb{R}^{(M-N) \times (M-N)}$. Note that the standard basis vectors of \mathbb{R}^N , namely, $\mathbf{e}^{(1)} = [1 \cdots 0]^\top, \dots, \mathbf{e}^{(N)} = [0 \cdots 1]^\top \in \mathbb{R}^N$, are the eigenvectors of $\frac{1}{2\lambda_2\sigma^2} \mathbf{S}(\mathbf{z})^{-2}$ and corresponding eigenvalues are $\frac{1}{2\lambda_2\sigma^2 s_1(\mathbf{z})^2} < \dots < \frac{1}{2\lambda_2\sigma^2 s_N(\mathbf{z})^2}$. According to the expression (40), the condition (36) means $\mathbf{C}^\top = \mathbf{O}_{(M-N) \times N}$, and thus $\mathbf{C}^\top \mathbf{e}^{(j)} = \mathbf{0}$ for all $j = 1, \dots, N$ in this situation. Note also that the eigenvalues of $\mathbf{U}(\mathbf{z})^\top \mathbf{A}(\mathbf{x}) \mathbf{U}(\mathbf{z})$ coincide with those of $\mathbf{A}(\mathbf{x})$. Therefore, if Eq. (36) holds, we have

$$\alpha_1(\mathbf{A}(\mathbf{x})) = \frac{1}{2\lambda_2\sigma^2 s_1(\mathbf{z})^2}, \dots, \alpha_N(\mathbf{A}(\mathbf{x})) = \frac{1}{2\lambda_2\sigma^2 s_N(\mathbf{z})^2}, \quad (41)$$

due to the inclusion principle (Horn & Johnson, 2013). Eq. (37) follows from Eq. (41).

As mentioned before, when the metric function is square of L2 norm, $\mathbf{A}(\mathbf{x})$ is the identity matrix \mathbf{I}_M . In this case, Eq. (36) holds and we have $J_{sv}(\mathbf{z}) = (1/2\lambda_1\sigma^2)^{N/2}$ †.

* Consider covering the original region in \mathbb{R}^N by infinitesimal hypercubes.

† This can also be directly confirmed by taking determinants of Eq. (35) after substituting $\mathbf{A}(\mathbf{x}) = \mathbf{I}_M$.

C. Effect of $h(x)$

In this section, the effects of $h(d)$ is discussed. By training the encoder and the decoder to be exact inverse functions of each other regarding the input data, the mapping becomes much rigidly isometric. Actually, for this purpose, it is important to choose $h(d)$ appropriately depending on metric function.

In this appendix we evaluate the behaviors of encoder and decoder in a one dimensional case using simple parametric linear encoder and decoder. Let x be a one dimensional data with the normal distribution:

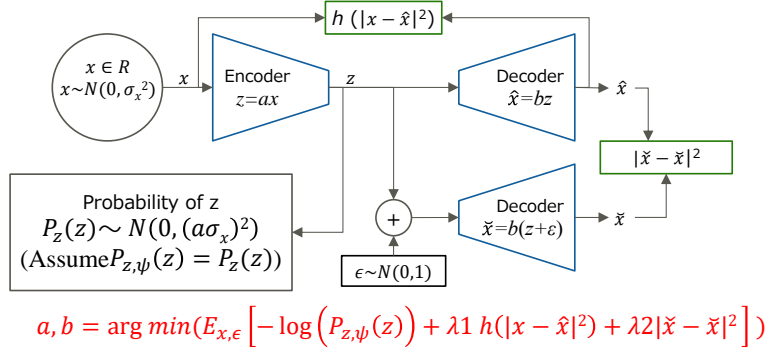


Figure 9. Simple encoder/decoder model to evaluate $h(d)$

$$\begin{aligned} x &\in \mathbb{R}, \\ x &\sim \mathcal{N}(0, \sigma_x^2). \end{aligned}$$

Let z be a one dimensional latent variable. Following two linear encoder and decoder are provided with parameter a and b :

$$\begin{aligned} z &= ax, \\ \hat{x} &= bz. \end{aligned}$$

Due to the above relationship, we have

$$P_z(z) = \mathcal{N}(0, (a\sigma_x)^2). \quad (42)$$

Here, square error is used as metrics function $D(x, y)$. The distribution of noise ϵ added to latent variable z is set to $N(0, 1)$. Then \check{x} is derived by decoding $z + \epsilon$ as:

$$\begin{aligned} D(x, y) &= |x - y|^2, \\ \epsilon &\sim \mathcal{N}(0, 1), \\ \check{x} &= b(z + \epsilon). \end{aligned}$$

For simplicity, we assume parametric PDF $P_{z, \psi}(z)$ is equal to the real PDF $P(z)$. Because the distribution of latent variable z follows $N(0, (a\sigma_x)^2)$, the entropy of z can be expressed as follows:

$$\begin{aligned} H(z) &= \int -P_z(z) \log(P_z(z)) dz \\ &= \log(a) + \log(\sigma_x \sqrt{2\pi e}). \end{aligned} \quad (43)$$

Using these notations, Eqs. (7) and (9) can be expressed as follows:

$$\begin{aligned} L &= E_{x \sim \mathcal{N}(0, \sigma_x^2), \epsilon \sim \mathcal{N}(0, 1)} [-\log P_z(z) + \lambda_1 h(|x - \hat{x}|^2) + \lambda_2 |\hat{x} - \check{x}|^2] \\ &= \log(a) + \log(\sigma_x \sqrt{2\pi e}) + \lambda_1 E_{x \sim \mathcal{N}(0, \sigma_x^2)} [h(|x - \hat{x}|^2)] + \lambda_2 b^2. \end{aligned} \quad (44)$$

At first, the case of $h(d) = d$ is examined. By applying $h(d) = d$, Eq. (44) can be expanded as follows:

$$L = \log(a) + \log(\sigma_x \sqrt{2\pi e}) + \lambda_1 (ab - 1)^2 \sigma_x^2 + \lambda_2 b^2. \quad (45)$$

By solving $\frac{\partial L}{\partial a} = 0$ and $\frac{\partial L}{\partial b} = 0$, a and b are derived as follows:

$$\begin{aligned} ab &= \frac{\lambda_1 \sigma_x^2 + \sqrt{\lambda_1^2 \sigma_x^4 - 2\lambda_1 \sigma_x^2}}{2\lambda_1 \sigma_x^2}, \\ a &= \sqrt{2\lambda_2} \left(\frac{\lambda_1 \sigma_x^2 + \sqrt{\lambda_1^2 \sigma_x^4 - 2\lambda_1 \sigma_x^2}}{2\lambda_1 \sigma_x^2} \right), \\ b &= 1/\sqrt{2\lambda_2}. \end{aligned}$$

If $\lambda_1 \sigma_x^2 \gg 1$, these equations are approximated as next:

$$\begin{aligned} ab &\simeq \left(1 - \frac{1}{2\lambda_1 \sigma_x^2} \right), \\ a &= \sqrt{2\lambda_2} \left(1 - \frac{1}{2\lambda_1 \sigma_x^2} \right), \\ b &= 1/\sqrt{2\lambda_2}. \end{aligned}$$

Here, ab is not equal to 1. That is, decoder is not an inverse function of encoder. In this case, the scale of latent space becomes slightly bent in order to minimize entropy function. As a result, good fitting of parametric PDF $P_z(z) \simeq P_{z,\psi}(z)$ could be realized while proportional relationship $P_z(z) \propto P_x(x)$ is relaxed.

Next, the case of $h(d) = \log(d)$ is examined. By applying $h(d) = \log(d)$ and introducing a minute variable Δ , Eq. (44) can be expanded as follows:

$$L = \log(a) + \log(\sigma_x \sqrt{2\pi e}) + \lambda_1 \log\left((ab - 1)^2 + \Delta\right) + \lambda_2 b^2. \quad (46)$$

By solving $\frac{\partial Loss}{\partial a} = 0$ and $\frac{\partial Loss}{\partial b} = 0$ and setting $\Delta \rightarrow 0$, a and b are derived as follows:

$$\begin{aligned} ab &= 1, \\ a &= \sqrt{2\lambda_2}, \\ b &= 1/\sqrt{2\lambda_2} \end{aligned} \quad (47)$$

Here, ab is equal to 1 and decoder becomes an inverse function of encoder regardless of the variance σ_x^2 . In this case, good proportional relation $P_z(z) \propto P_x(x)$ could be realized regardless of the fitting $P_{z,\psi}(z)$ to $P_z(z)$.

Considering from these results, there could be a guideline to choose $h(d)$. If the parametric PDF $P_{z,\psi}(z)$ has enough ability to fit the real distribution $P_z(z)$, $h(d) = \log(d)$ could be better. If not, $h(d) = d$ could be an option.

D. Expansion of SSIM and BCE to Quadratic Forms

In this appendix, it is shown that SSIM and BCE can be approximated in quadratic forms with positive definite matrices except some constants.

Structural similarity (SSIM) (Wang et al., 2004) is widely used for picture quality metric since it is close to human subjective evaluation. In this appendix, we show $(1 - SSIM)$ can be approximated to a quadratic form such as Eq.(8).

Eq. (48) is a SSIM value between cropped pictures \mathbf{x} and \mathbf{y} with a $W \times W$ window:

$$SSIM_{W \times W}(\mathbf{x}, \mathbf{y}) = \frac{2\mu_x \mu_y}{\mu_x^2 + \mu_y^2} \frac{2\sigma_{xy}}{\sigma_x^2 + \sigma_y^2}. \quad (48)$$

In order to calculate SSIM index for entire pictures, this window is shifted in a whole picture and all of SSIM values are averaged. If $(1 - SSIM_{W \times W}(\mathbf{x}, \mathbf{y}))$ is expressed in quadratic form, the average for a picture $(1 - SSIM_{picture})$ can be also expressed in quadratic form.

Let $\Delta \mathbf{x}$ be a minute displacement of \mathbf{x} . Then SSIM between \mathbf{x} and $\mathbf{x} + \Delta \mathbf{x}$ can be expressed as follows:

$$SSIM_{W \times W}(\mathbf{x}, \mathbf{x} + \Delta \mathbf{x}) = 1 - \frac{\mu_{\Delta \mathbf{x}}^2}{2\mu_x^2} - \frac{\sigma_{\Delta \mathbf{x}}^2}{2\sigma_x^2} + O\left(\frac{|\Delta \mathbf{x}|}{|\mathbf{x}|}\right)^3 \quad (49)$$

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

$$\mu_{\Delta \mathbf{x}}^2 = \Delta \mathbf{x}^\top \mathbf{W}_m \Delta \mathbf{x}, \quad (50)$$

$$\sigma_{\Delta \mathbf{x}}^2 = \Delta \mathbf{x}^\top \mathbf{W}_v \Delta \mathbf{x}, \quad (51)$$

where

$$\mathbf{W}_m = \frac{1}{W^2} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}, \quad \mathbf{W}_v = \frac{1}{W^2} \begin{pmatrix} W-1 & -1 & \dots & -1 \\ -1 & W-1 & \dots & -1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \dots & W-1 \end{pmatrix}. \quad (52)$$

It should be noted that matrix \mathbf{W}_m is positive definite and matrix \mathbf{W}_v is positive semidefinite. As a result, $(1 - SSIM_{W \times W}(\mathbf{x}, \mathbf{y}))$ can be expressed in the following quadratic form with positive definite matrix:

$$1 - SSIM_{W \times W}(\mathbf{x}, \mathbf{x} + \Delta \mathbf{x}) \simeq \Delta \mathbf{x}^\top \left(\frac{1}{2\mu_x^2} \mathbf{W}_m + \frac{1}{2\sigma_x^2} \mathbf{W}_v \right) \Delta \mathbf{x}. \quad (53)$$

Binary cross entropy (BCE) is also a reconstruction loss function widely used in VAE (Kingma & Welling, 2014). BCE is defined as follows:

$$BCE(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^M (-x_i \log(y_i) - (1 - x_i) \log(1 - y_i)). \quad (54)$$

BCE can be also approximated by a quadratic form with positive definite matrix. Let $\Delta \mathbf{x}$ be a small displacement of \mathbf{x} and Δx_i be its i -th component. Then BCE between \mathbf{x} and $\mathbf{x} + \Delta \mathbf{x}$ can be expanded as follows:

$$\begin{aligned} BCE(\mathbf{x}, \mathbf{x} + \Delta \mathbf{x}) &= \sum_i (-x_i \log(x_i + \Delta x_i) - (1 - x_i) \log(1 - x_i - \Delta x_i)) \\ &= \sum_i \left(-x_i \log \left(x_i \left(1 + \frac{\Delta x_i}{x_i} \right) \right) - (1 - x_i) \log \left((1 - x_i) \left(1 - \frac{\Delta x_i}{1 - x_i} \right) \right) \right) \\ &= \sum_i \left(-x_i \log \left(1 + \frac{\Delta x_i}{x_i} \right) - (1 - x_i) \log \left(1 - \frac{\Delta x_i}{1 - x_i} \right) \right) \\ &\quad + \sum_i (-x_i \log(x_i) - (1 - x_i) \log(1 - x_i)). \end{aligned} \quad (55)$$

Here, the second term of the last equation is constant depending on \mathbf{x} . The first term of the last equation is further expanded as follows by using Maclaurin expansion of logarithm:

$$\begin{aligned} &\sum_i \left(-x_i \left(\frac{\Delta x_i}{x_i} - \frac{\Delta x_i^2}{2x_i^2} \right) - (1 - x_i) \left(-\frac{\Delta x_i}{1 - x_i} - \frac{\Delta x_i^2}{2(1 - x_i)^2} \right) + O(\Delta x_i^3) \right) \\ &= \sum_i \left(\frac{1}{2} \left(\frac{1}{x_i} + \frac{1}{1 - x_i} \right) \Delta x_i^2 + O(\Delta x_i^3) \right). \end{aligned} \quad (56)$$

Then, let a matrix $\mathbf{A}(\mathbf{x})$ be defined as follows:

$$\mathbf{A}(\mathbf{x}) = \begin{pmatrix} \frac{1}{2} \left(\frac{1}{x_1} + \frac{1}{1-x_1} \right) & 0 & \dots \\ 0 & \frac{1}{2} \left(\frac{1}{x_2} + \frac{1}{1-x_2} \right) & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (57)$$

Obviously $\mathbf{A}(\mathbf{x})$ is a positive definite matrix. As a result, BCE between \mathbf{x} and $\mathbf{x} + \Delta\mathbf{x}$ can be approximated by a quadratic form with \mathbf{x} depending constant offset as follows:

$$BCE(\mathbf{x}, \mathbf{x} + \Delta\mathbf{x}) \simeq \Delta\mathbf{x}^\top \mathbf{A}(\mathbf{x})\Delta\mathbf{x} + \sum_i (-x_i \log(x_i) - (1 - x_i) \log(1 - x_i)). \quad (58)$$

Note that BCE is typically used for binary data. In this case, the second term in Eq. (58) is always 0.

E. “Continuous PCA” Feature of Isometric Embedding for Riemannian Manifold

In this section, we explain that the isometric embedding realized by RaDOGAGA has a continuous PCA feature when the following factorized probability density model is used:

$$P_{\mathbf{z},\psi}(\mathbf{z}) = \prod_{i=1}^N P_{z_i,\psi}(z_i). \quad (59)$$

Here, our definition of “continuous PCA” is the following. 1) Mutual information between latent variables are minimum and likely to be uncorrelated to each other: 2) Energy of latent space is concentrated to several principal components, and the importance of each component can be determined: 3) These features are held for all subspace of a manifold and subspace is continuously connected.

Next we explain the reason why these feature is acquired. As explained in Appendix A, all column vectors of Jacobian matrix of decoder from latent space to data space have the same norm and all combinations of pairwise vectors are orthogonal. In other words, when constant value is multiplied, the resulting vectors are orthonormal. Because encoder is a inverse function of decoder ideally, each row vector of encoder’s Jacobian matrix should be the same as column vector of decoder under the ideal condition. Here, $f_{ortho,\theta}(\mathbf{x})$ and $g_{ortho,\phi}(\mathbf{z}_\theta)$ are defined as encoder and decoder with these feature. Because the latent variables depend on encoder parameter θ , latent variable is described as $\mathbf{z}_\theta = f_{ortho,\theta}(\mathbf{x})$, and its PDF is defined as $P_{\mathbf{z},\theta}(\mathbf{z}_\theta)$. PDFs of latent space and data space have the following relation where $J_{sv}(\mathbf{z}_\theta)$ is the product of the singular values of $\mathbf{J}(\mathbf{z}_\theta)$ which is a Jacobian matrix between two spaces as explained in Section 4.2 and Appendix B.

$$P_{\mathbf{z},\theta}(\mathbf{z}_\theta) = J_{sv}(\mathbf{z}_\theta) P_{\mathbf{x}}(\mathbf{x}) \propto \left(\prod_{j=1}^N \alpha_j(\mathbf{A}(\mathbf{x})) \right)^{-\frac{1}{2}} P_{\mathbf{x}}(\mathbf{x}). \quad (60)$$

As described before, $P_{\mathbf{z},\psi}(\mathbf{z})$ is a parametric PDF of the latent space to be optimized with parameter ψ .

By applying the result of Eqs. (24) and (31), Eq. (7) can be transformed as Eq. (61) where $\hat{\mathbf{x}} = g_{ortho,\phi}(f_{ortho,\theta}(\mathbf{x}))$.

$$\begin{aligned} L_{ortho} &= -\log(P_{\mathbf{z},\psi}(\mathbf{z}_\theta)) + \lambda_1 h(D(\mathbf{x}, \hat{\mathbf{x}})) + N/2. \\ \text{s.t. } \left(\frac{\partial g_{ortho,\phi}(\mathbf{z}_\theta)}{\partial z_{\theta_i}} \right)^\top \mathbf{A}(\mathbf{x}) \left(\frac{\partial g_{ortho,\phi}(\mathbf{z}_\theta)}{\partial z_{\theta_j}} \right) &= \frac{1}{2\lambda\sigma^2} \delta_{ij}. \end{aligned} \quad (61)$$

Here, the third term of the right side is constant, this term can be removed from the cost function as follows:

$$L'_{ortho} = -\log(P_{\mathbf{z},\psi}(\mathbf{z}_\theta)) + \lambda_1 h(D(\mathbf{x}, \hat{\mathbf{x}})). \quad (62)$$

Then the parameters of network and PDF are obtained according to the following equation:

$$\theta, \phi, \psi = \arg \min_{\theta, \phi, \psi} (E_{\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})} [L'_{ortho}]). \quad (63)$$

$E_{\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})} [L'_{ortho}]$ in Eq. (63) can be transformed as the next:

$$\begin{aligned} E_{\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})} [L'_{ortho}] &= \int P_{\mathbf{x}}(\mathbf{x}) (-\log(P_{\mathbf{z},\psi}(\mathbf{z}_\theta)) + \lambda_1 h(D(\mathbf{x}, \hat{\mathbf{x}}))) d\mathbf{x} \\ &= \int (P_{\mathbf{z},\theta}(\mathbf{z}_\theta) J_{sv}(\mathbf{z}_\theta)^{-1}) (-\log(P_{\mathbf{z},\psi}(\mathbf{z}_\theta))) J_{sv}(\mathbf{z}_\theta) dz_\theta + \lambda_1 \int P_{\mathbf{x}}(\mathbf{x}) h(D(\mathbf{x}, \hat{\mathbf{x}})) d\mathbf{x}. \end{aligned} \quad (64)$$

At first, the first term of the third formula in Eq.(64) is examined. Let $d\mathbf{z}_{\theta/i}$ be a differential of $(N - 1)$ dimensional latent variables where i -th axis $z_{\theta i}$ is removed from the latent variable \mathbf{z}_{θ} . Then a marginal distribution of $z_{\theta i}$ can be derived from the next equation:

$$P_{z,\theta i}(z_{\theta i}) = \int P_{\mathbf{z},\theta}(\mathbf{z}_{\theta}) d\mathbf{z}_{\theta/i}. \quad (65)$$

By using Eqs.(59) and (65), the first term of the third formula in Eq. (64) can be expanded as:

$$\begin{aligned} \int P_{\mathbf{z},\theta}(\mathbf{z}_{\theta}) (-\log(P_{\mathbf{z},\psi}(\mathbf{z}_{\theta}))) d\mathbf{z}_{\theta} &= \int P_{\mathbf{z},\theta}(\mathbf{z}_{\theta}) \left(-\log \left(\frac{\prod_{i=1}^N P_{z_i,\psi}(z_{\theta i})}{\prod_{i=1}^N P_{z,\theta i}(z_{\theta i})} \right) \right) d\mathbf{z}_{\theta} \\ &\quad + \int P_{\mathbf{z},\theta}(\mathbf{z}_{\theta}) \left(-\log \left(\prod_{i=1}^N P_{z,\theta i}(z_{\theta i}) \right) \right) d\mathbf{z}_{\theta} \\ &= \sum_{i=1}^N \int \left(\int P_{\mathbf{z},\theta}(\mathbf{z}_{\theta}) d\mathbf{z}_{\theta/i} \right) \left(-\log \left(\frac{P_{z_i,\psi}(z_{\theta i})}{P_{z,\theta i}(z_{\theta i})} \right) \right) dz_{\theta i} \\ &\quad + \sum_{i=1}^N \int \left(\int P_{\mathbf{z},\theta}(\mathbf{z}_{\theta}) d\mathbf{z}_{\theta/i} \right) (-\log(P_{z,\theta i}(z_{\theta i}))) dz_{\theta i} \\ &= \sum_{i=1}^N D_{KL}(P_{z,\theta i}(z_{\theta i}) \| P_{z_i,\psi}(z_{\theta i})) + \sum_{i=1}^N H(z_{\theta i}). \end{aligned} \quad (66)$$

Here $H(X)$ denotes the entropy of a variable X . The first term of the third formula is KL-divergence between marginal probability $P_{z,\theta i}(z_{\theta i})$ and factorized parametric probability $P_{z_i,\psi}(z_{\theta i})$. The second term of the third formula can be further transformed using mutual information between latent variables $I(\mathbf{z}_{\theta})$ and equation (60).

$$\begin{aligned} \sum_{i=1}^N H(z_{\theta i}) &= H(\mathbf{z}_{\theta}) + I(\mathbf{z}_{\theta}) \simeq - \int J_{sv}(\mathbf{z}_{\theta}) P_{\mathbf{x}}(\mathbf{x}) \log(J_{sv}(\mathbf{z}_{\theta}) P_{\mathbf{x}}(\mathbf{x})) J_{sv}(\mathbf{z}_{\theta})^{-1} d\mathbf{x} + I(\mathbf{z}_{\theta}) \\ &= H(\mathbf{x}) - \int P_{\mathbf{x}}(\mathbf{x}) \log \left(\left(\frac{1}{2\lambda_2\sigma^2} \right)^{\frac{N}{2}} \left(\prod_{j=1}^N \alpha_j(\mathbf{A}(\mathbf{x})) \right)^{-\frac{1}{2}} \right) d\mathbf{x} + I(\mathbf{z}_{\theta}) \end{aligned} \quad (67)$$

At second, the second term of the third formula in Eq. (64) is examined. When \mathbf{x} and $\hat{\mathbf{x}}$ are close, the following equation holds.

$$D(\mathbf{x}, \hat{\mathbf{x}}) \simeq (\mathbf{x} - \hat{\mathbf{x}})^{\top} \mathbf{A}(\mathbf{x})(\mathbf{x} - \hat{\mathbf{x}}). \quad (68)$$

Note that with given distribution $\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})$, the first and the second term in the right side of Eq. (67) are fixed value. Therefore, by using these expansions, Eq.(64) can be expressed as:

$$\begin{aligned} E_{\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})} [L'_{ortho}] &\simeq \sum_{i=1}^N D_{KL}(P_{z,\theta i}(z_{\theta i}) \| P_{z_i,\psi}(z_{\theta i})) \\ &\quad + I(\mathbf{z}_{\theta}) + E_{\mathbf{x}} [(\mathbf{x} - \hat{\mathbf{x}})^{\top} \mathbf{A}(\mathbf{x})(\mathbf{x} - \hat{\mathbf{x}})] + \text{Const}. \end{aligned} \quad (69)$$

Here, the real space \mathbb{R}^M is divided into a plurality of small subspace partitioning $\Omega_{\mathbf{x}_1}, \Omega_{\mathbf{x}_2}, \dots$. Note that \mathbb{R}^M is an inner product space endowed with metric tensor $\mathbf{A}(\mathbf{x})$. Let $\Omega_{z_1}, \Omega_{z_2}, \dots$ be the division space of the latent space $\mathbf{z} \in \mathbb{R}^N$ corresponding to $\Omega_{\mathbf{x}}$.

Then Eq. (69) can be rewritten as:

$$\begin{aligned} E_{\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})} [L'_{ortho}] &\simeq \sum_{i=1}^N D_{KL}(P_{z,\theta i}(z_{\theta i}) \| P_{z_i,\psi}(z_{\theta i})) \\ &\quad + \sum_k (I(\mathbf{z}_{\theta} \in \Omega_{z_{\theta k}}) + E_{\mathbf{x} \in \Omega_{\mathbf{x}_k}} [(\mathbf{x} - \hat{\mathbf{x}})^{\top} \mathbf{A}(\mathbf{x})(\mathbf{x} - \hat{\mathbf{x}})]) + \text{Const}. \end{aligned} \quad (70)$$

For each subspace partitioning, Jacobian matrix for the transformation from Ωx_k to $\Omega z_{\theta k}$ forms constantly scaled orthonormal system with respect to $A(x)$. According to Karhunen-Loève Theory (Rao & Yip, 2000), the orthonormal basis which minimize both mutual information and reconstruction error leads to be Karhunen-Loève transform (KLT). It is noted that the basis of KLT is equivalent to PCA orthonormal basis.

As a result, when Eq. (70) is minimized, Jacobi matrix from Ωx_k to $\Omega z_{\theta k}$ for each subspace partitioning should be KLT/PCA. Consequently, the same feature as PCA will be realized such as the determination of principal components etc.

From these considerations, we conclude that RaDOGAGA has a “continuous PCA” feature. This is experimentally shown in Section 5.1 and Appendix F.6.

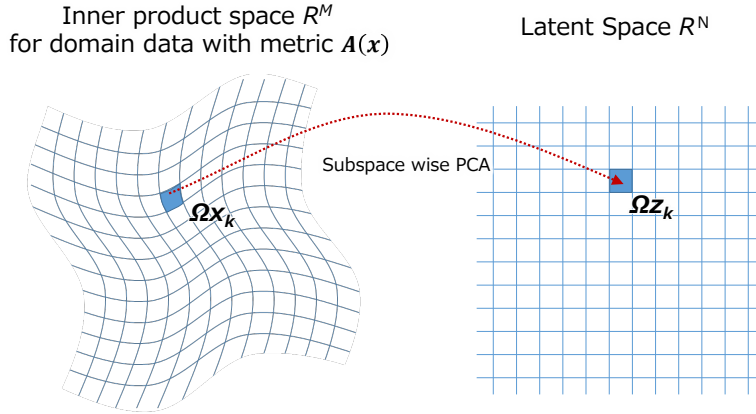


Figure 10. Continuous KLT (PCA) Mapping from input domain to latent space. For all small subspace partitioning Ωx_k domain space (which is inner product space with metric tensor $A(x)$), mapping from Ωx_k to $\Omega z_{\theta k}$ can be regarded as PCA

F. Detail and Expansion Result of Experiment in Section 5.1

In this section, we will provide further detail and a result of the complementary experiment regarding section 5.1.

F.1. GDN Activation

GDN activation function (Ballé et al., 2016) is known to be suitable for image compression. For implementation, we use a TensorFlow library^{||}.

F.2. Other Training Information

The batch size is 64. The iteration number is 500,000. We use NVIDIA Tesla V100 (SXM2). For RaDOGAGA, since our implementation was done based on the source code for image compression, entropy rate is calculated as $-\log((P_{z,\psi}(z))/(M \log 2))$, meaning bit per pixel. In addition, for RaDOGAGA, the second term of Eq. (7) is always MSE in this experiment. This is because we found that the training with $1 - SSIM$ as the reconstruction loss is likely to diverge at the beginning step of the training. Therefore, we tried to start training with MSE and then fine-tuned with $1 - SSIM$. Eventually, the result is almost the same as the case without finetuning. Therefore, to simplify the training process, we do not usually finetune.

F.3. Generation of v_z and w_z

To evaluate the isometricity of the mapping, it is necessary to prepare random tangent vector v_z and w_z with a scattered interior angle. We generate two different tangent vectors $v_z = \{v_{z1}, v_{z2}, \dots, v_{zn}\}$ and $w_z = \{w_{z1}, w_{z2}, \dots, w_{zn}\}$ in the following manner. First, we prepare $v'_z \in \mathbb{R}^N$ as $\{1.0, 0.0, \dots, 0.0\}$. Then, we sample $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_{n-1}\}$ ($\alpha_{1 \dots n-2} \sim U(0, \pi), \alpha_{n-1} \sim U(0, 2\pi)$) to set w' as the conversion of polar coordinate $\{r, \alpha\} \in \mathbb{R}^N$ to rectangular

^{||}https://github.com/tensorflow/compression/tree/master/docs/api_docs/python/tfc

coordinates, where $r = 1$. Thus, the distribution of interior angle of \mathbf{v}'_z and \mathbf{w}'_z also obey $\alpha_1 \sim U(0, \pi)$. Next, we randomly rotate the plane \mathbb{R}^N in which interior angle of \mathbf{v}'_z and \mathbf{w}'_z is α_1 in the following way (Teoh, 2005) and obtain \mathbf{v}_z and \mathbf{w}_z .

$$\boldsymbol{\rho} = -\frac{\cos \alpha_1}{\sin \alpha_1} \mathbf{v}'_z + \frac{1}{\sin \alpha_1} \mathbf{w}'_z, \quad \boldsymbol{\tau} = \mathbf{v}'_z,$$

then,

$$\begin{pmatrix} \mathbf{v}_z \\ \mathbf{w}_z \end{pmatrix} = \begin{bmatrix} -\sin \omega & \cos \omega \\ \cos \omega \sin \alpha_1 - \sin \omega \cos \alpha_1 & \sin \omega \sin \alpha_1 + \cos \omega \cos \alpha_1 \end{bmatrix} \begin{pmatrix} \boldsymbol{\rho} \\ \boldsymbol{\tau} \end{pmatrix}, \quad (71)$$

where $\omega \sim U(0, 2\pi)$ is the rotation angle of the plane. Note that since this is the rotation of the plane, the interior angle between \mathbf{v} and \mathbf{w} is kept to α_1 . Finally, we normalize the norm of \mathbf{v}_z and \mathbf{w}_z to be 0.01.

F.4. Experiment with MNIST Dataset and BCE

Besides of the experiment in main paper, we conducted an experiment with MNIST dataset (LeCun et al., 1998)^{††} which contains handwritten digits binary images with the image size of 28×28 . We use 60,000 samples in the training split. The metric function is *BCE*, where $\mathbf{A}(\mathbf{x})$ is approximated as Eq. (57). Autoencoder consists of FC layers with sizes of 1000, 1000, 128, 1000, and 1000. We attach *softplus* as activation function except for the last of the encoder and the decoder. In this experiment, we modify the form of the cost function of beta-VAE as

$$L = -L_{kl} + \lambda_1 h(D(\mathbf{x}, \hat{\mathbf{x}})) + \lambda_2 D(\hat{\mathbf{x}}, \check{\mathbf{x}}), \quad (72)$$

where $\hat{\mathbf{x}}$ is the output of the decoder without noise, and $\check{\mathbf{x}}$ is the output of the decoder with the noise of reparameterization trick. We set (λ_1, λ_2) as (10, 1) for beta-VAE and (0.01, 0.01) for RaDOGAGA. Optimization is done with Adam optimizer with learning rate 1×10^{-4} for beta-VAE 1×10^{-5} for RaDOGAGA. The batch size is 256 and the training iteration is 30,000. These parameters are determined to make the $PSNR = 20 \log_{10} \left(\frac{MAX_{\mathbf{x}}^2}{MSE} \right)$, where $MAX_{\mathbf{x}} = 255$, between input and reconstruction image approximately 25 dB.

Figure 11 depicts the result. We can observe that map of RaDOGAGA is isometric as well even for the case the metric function is *BCE*. Consequently, even if the metric function is complicated one, the impact of the latent variable on the metric function is tractable. We expect this feature promotes further improving of metric learning, data interpolation, and so on.

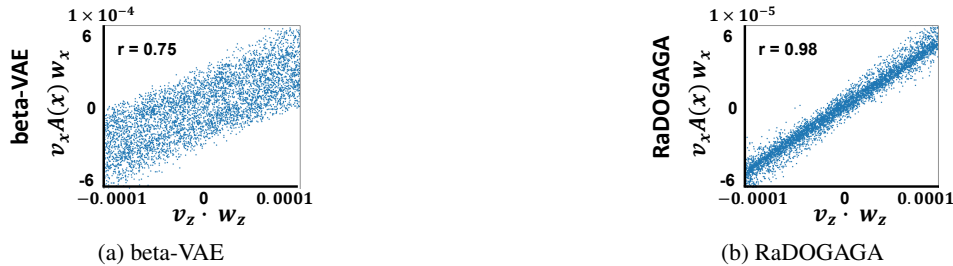


Figure 11. Plot of $\mathbf{v}_z \cdot \mathbf{w}_z$ (horizontal axis) and $\mathbf{v}_x^\top \mathbf{A}(\mathbf{x}) \mathbf{w}_x$ (vertical axis). In beta-VAE (left), the correlation is weak while in our method (right) we can observe proportionality.

F.5. Isometricity of Encoder Side

In Section 5.1, we showed the isometricity of decoder side because it is common to analyse the behavior of latent variables by observing the decoder output such as latent traverse. We also clarify that the embedding by encoder f keep isometric. Given two tangent vector \mathbf{v}_x and \mathbf{w}_x , $\mathbf{v}_x^\top \mathbf{A}(\mathbf{x}) \mathbf{w}_x$ is compared to $d\mathbf{f}(\mathbf{v}_x) \cdot d\mathbf{f}(\mathbf{w}_x)$. $d\mathbf{f}(\mathbf{w}_x)$ is also approximated by $f(\mathbf{x} + \mathbf{v}_x) - f(\mathbf{x})$. As Fig. 12 shows, the embedding to the latent space is isometric. Consequently, it is experimentally supported that our method enables to embed data in Euclidean space isometrically. The result of the same experiment for the case of the metric is $1 - SSIM$ is provided in Appendix F.

^{††}<http://yann.lecun.com/exdb/mnist/>



Figure 12. $v_x^\top A(x)w_x$ vs $df(v_x) \cdot df(w_x)$. The mapping by encoder is also isometric.

F.6. Additional Latent Traverse

In Section 5.1, the latent traverse for variables with the top 9 variances was provided. To further clarify whether the variance is corresponding to visual impact, the latent traverse of RaDOGAGA for $z_0, z_1, z_2, z_{20}, z_{21}, z_{22}, z_{200}, z_{201},$ and z_{202} are shown in Fig. 13. Apparently, a latent traverse with a larger σ makes a bigger difference in the image. When the σ^2 gets close to 0, there is almost no visual difference. Accordingly, the behavior as continuous PCA is clarified throughout the entire variables.

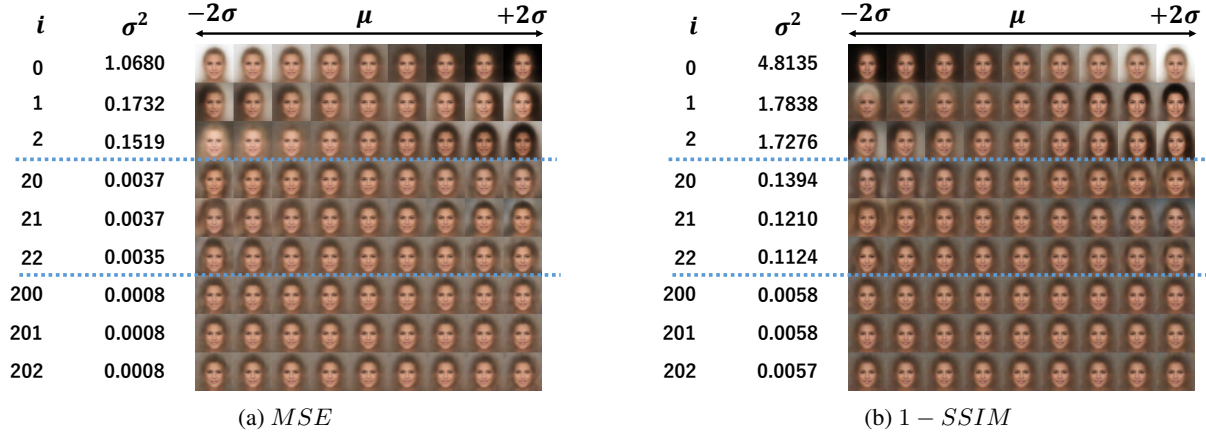


Figure 13. Latent space traversal of z . For the top-3 variables, pictures look significantly different. In the middle range (z_{20}, z_{21}, z_{22}), the difference is smaller than the upper three but still observable. For the bottom three, there is almost no difference.

G. Detail of the Experiment in Section 5.3

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

G.1. Datasets

We describe the detail of following four public datasets:

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

Thyroid (Dua & Graff, 2019) This dataset contains 3,772 data sample with 6-dimensional feature from patients and can be divided in three classes: normal (not hypothyroid), hyperfunction, and subnormal functioning. We treat 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 containing 452 data sample 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, and the others are treated as normal.

KDDCUP-Rev (Dua & Graff, 2019) To treat “normal” instances as majority in the KDDCUP dataset, we keep all “normal” instances and randomly pick up “attack” instances so that they compose 20% of the dataset. In the end, the number of instance is 121,597.

Data is max-min normalized toward dimension through the entire dataset.

G.2. Hyperparameter and Training Detail

Hyperparameter for RaDOGAGA is described in Table 2. First and second column is number of neurons. (λ_1, λ_2) is determined experimentally. For DAGMM, the number of neuron is the same as Table 2. We set (λ_1, λ_2) as $(0.1, 0.005)$ referring Zong et al. (2018) except for Thyroid. Only for Thyroid, (λ_1, λ_2) is $(0.1, 0.0001)$ since $(0.1, 0.005)$ does not work well with our implementation. Optimization is done by Adam optimizer with learning rate 1×10^{-4} for all dataset. The batch size is 1024 for all dataset. The epoch number is 100, 20000, 10000, and 100 respectively. We save and the test models by every 1/10 epochs and early stop is applied. For this experiment, we use GeForce GTX 1080.

Table 2. Hyper parameter for RaDOGAGA

Dataset	Autoencoder	EN	$\lambda_1(d)$	$\lambda_2(d)$	$\lambda_1(\log(d))$	$\lambda_2(\log(d))$
KDDCup99	60, 30, 8, 30, 60	10, 4	100	1000	10	100
Thyroid	30, 24, 6, 24, 30	10, 2	10000	1000	100	1000
Arrhythmia	10, 4, 10	10, 2	1000	1000	1000	100
KDDCup-rev	60, 30, 8, 30, 60	10, 2	100	100	100	100

G.3. Experiment with different network size

In addition to experiment in main page, we also conducted experiment with same network size as in Zong et al. (2018) with parameters in Table 3

Table 3. Hyper parameter for RaDOGAGA(same network size as in Zong et al. (2018))

Dataset	Autoencoder	EN	$\lambda_1(d)$	$\lambda_2(d)$	$\lambda_1(\log(d))$	$\lambda_2(\log(d))$
KDDCup99	60, 30, 1, 30, 60	10, 4	100	100	100	1000
Thyroid	12, 4, 1, 4, 12	10, 2	1000	10000	100	10000
Arrhythmia	10, 2, 10	10, 2	1000	100	1000	100
KDDCup-rev	60, 30, 1, 30, 60	10, 2	100	100	100	1000

Now, we provide results of setting in Table 3. In Table 4, RaDOGAGA- and DAGMM- are results of them and DAGMM is result cited from Zong et al. (2018). Even with this network size, our method has boost from baseline in all dataset.

Table 4. Average and standard deviations (in brackets) of Precision, Recall and F1

Dataset	Methods	Precision	Recall	F1
KDDCup	DAGMM	0.9297	0.9442	0.9369
	DAGMM-	0.9338 (0.0051)	0.9484 (0.0052)	0.9410 (0.0051)
	RaDOGAGA-(L2)	0.9455 (0.0016)	0.9608 (0.0018)	0.9531 (0.0017)
	RaDOGAGA-(log)	0.9370 (0.0024)	0.9517 (0.0025)	0.9443 (0.0024)
Thyroid	DAGMM	0.4766	0.4834	0.4782
	DAGMM-	0.4635 (0.1054)	0.4837 (0.1100)	0.4734 (0.1076)
	RaDOGAGA-(L2)	0.5729 (0.0449)	0.5978 (0.0469)	0.5851 (0.0459)
	RaDOGAGA-(log)	0.5729 (0.0398)	0.5978 (0.0415)	0.5851 (0.0406)
Arrythmia	DAGMM	0.4909	0.5078	0.4983
	DAGMM-	0.4721 (0.0451)	0.4864 (0.0464)	0.4791 (0.0457)
	RaDOGAGA-(L2)	0.4897 (0.0477)	0.5045 (0.0491)	0.4970 (0.0484)
	RaDOGAGA-(log)	0.5044 (0.0364)	0.5197 (0.0375)	0.5119 (0.0369)
KDDCup-rev	DAGMM	0.937	0.939	0.938
	DAGMM-	0.9491 (0.0163)	0.9498 (0.0158)	0.9494 (0.0160)
	RaDOGAGA-(L2)	0.9761 (0.0057)	0.9761 (0.0056)	0.9761 (0.0057)
	RaDOGAGA-(log)	0.9791 (0.0036)	0.9799 (0.0035)	0.9795 (0.0036)