

# RATE-DISTORTION OPTIMIZATION GUIDED AUTOENCODER FOR GENERATIVE APPROACH

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In the generative model approach of machine learning, it is essential to acquire an accurate probabilistic model and compress the dimension of data for easy treatment. However, in the conventional deep-autoencoder based generative model such as VAE, the probability of the real space cannot be obtained correctly from that of in the latent space, because the scaling between both spaces is not controlled. This has also been an obstacle to quantifying the impact of the variation of latent variables on data. In this paper, we propose a method to learn parametric probability distribution and autoencoder simultaneously based on Rate-Distortion Optimization to support scaling control. It is proved theoretically and experimentally that (i) the probability distribution of the latent space obtained by this model is proportional to the probability distribution of the real space because Jacobian between two spaces is constant; (ii) our model behaves as non-linear PCA, which enables to evaluate the influence of latent variables on data. Furthermore, to verify the usefulness on the practical application, we evaluate its performance in unsupervised anomaly detection and it outperforms current state-of-the-art methods.

## 1 INTRODUCTION

Capturing the inherent features of a dataset from high-dimensional and complex data is an essential issue in machine learning. Generative model approach learns the probability distribution of data, aiming at data generation by probabilistic sampling, unsupervised/weakly supervised learning, and acquiring meta-prior (general assumptions about how data can be summarized naturally, such as disentangle, clustering, and hierarchical structure (Y. Bengio & Vincent, 2013; Michael Tschannen & Lucic, 2019)). It is generally difficult to directly estimate a probability density function(PDF)  $Px(x)$  of real data  $x$ . Accordingly, one promising approach is to map to the latent space  $z$  with reduced dimension and capture PDF  $Pz(z)$ . In recent years, deep autoencoder based methods have made it possible to compress dimensions and derive latent variables. While there is remarkable progress in these areas (van den Oord et al.; Kingma et al., 2014; Jiang et al., 2016), the relation between  $x$  and  $z$  in the current deep generative models is still not clear.

VAE (P.Kingma & Welling, 2014) is one of the most successful generative models for capturing latent representation. In VAE, lower bound of log-likelihood of  $Px(x)$  is introduced as ELBO. Then latent variable is obtained by maximizing ELBO. Some previous works reported that there are various limitations in the origin of maximizing ELBO (Alemi et al., 2018; Zhao et al., 2019). Moreover, many previous works did not care about Jacobian between two spaces, despite the fact that the ratio between  $Pz(z)$  and  $Px(x)$  is equal to the Jacobian. Even in models that provide more flexible estimation (M.Johnson, 2016; Liao et al., 2018; Zong et al., 2018), this point is overlooked.

Here, when we turn our sight toward acquiring meta-prior, it is straightforward to evaluate the quantitative influence of each latent variable on the distance between data  $x_1$  and  $x_2$ . To do so, the scale of the latent variable should be appropriately controlled so that the changes in latent variables is proportional to the changes of distance in data space. In addition, this scaling should be adjusted according to the definition of the distance metric. For instance, with respect to image quality metrics, different meta-prior would be derived from MSE and SSIM.

To deal with this, we propose RaDOGAGA (Rate-Distortion Optimization Guided Autoencoder for Generative Approach), which simultaneously learns parametric probability distribution and auto-

encoder, based on the rate-distortion optimization (RDO). In this paper, we show the effect of RaDOGAGA in the following steps.

- (1) We prove that RaDOGAGA has the following property theoretically and experimentally.
  - Jacobi matrix between real space and latent space leads to be constant-scaled orthonormal. So the response of the minute change of  $z$  to the real space data  $x$  is constant at any  $z$ .
  - Because of constant Jacobian (or pseudo-Jacobian),  $Px(x)$  and  $Pz(z)$  are almost proportional. Therefore,  $Px(x)$  can be estimated, by directly maximizing log-likelihood of parametric PDF  $Pz_\psi(z)$  in reduced-dimensional space, without considering ELBO.
  - When univariate independent distribution is used to estimate  $Pz(z)$  parametrically, it behaves as "continuous PCA" where energy is concentrated on several principal components.
- (2) Thanks to this property, RaDOGAGA achieve the state-of-the-art performance in anomaly detection task with four public datasets, where probability density estimation is important.
- (3) We show that our approach can directly evaluate how the  $z$  impact on the distance metric in real space. This feature is promising to further interpretation of latent variables.

## 2 RELATED WORK

**Flow based model:** Flow based generative models generates astonishing quality of image (Kingma & Dhariwal, 2018; Dinh et al., 2014). Flow mechanism explicitly takes Jacobian of  $x$  and  $z$  into account. The transformation function  $z = f(x)$  is learned, calculating and preserving Jacobian of  $x$  and  $z$ . Unlike ordinary autoencoder, which reverse  $z$  to  $x$  with function  $g()$  different from  $f()$ , inverse function transforms  $z$  as  $x = f^{-1}(z)$ . Since the model preserves Jacobian,  $Pz(x)$  can be estimated by maximizing log likelihood of  $Pz(z)$  without considering ELBO. Although, in this approach,  $f()$  need to be bijection. Because of this limitation, it is difficult to fully utilize the flexibility of neural networks.

**Interpretation of latent variables:** While it is expected to acquire meta-prior by deep autoencoder, interpreting latent variables is still challenging. In recent years, research aiming to acquire disentangled latent variables, which encourages them to be independent, is flourishing (Lopez et al., 2018; Chen et al., 2018; I. Higgins & Lerchner, 2018; Chen et al., 2016). With these methods, qualitative effects for disentanglement can be seen. For example, when a certain latent variable is displaced, image changes corresponding to specific attributes (size, color, etc.). Some works also propose quantitative metrics for meta-prior. In beta-VAE (I. Higgins & Lerchner, 2017), the metric evaluates the independence of latent variables by solving the classification task. But, the metrics do not evaluate the effect of latent variable to the distance between data directly, in spite of the fact that model is trained based on the distance.

**Image compression with rate-distortion optimization:** Rate-distortion optimization (RDO) minimizes the information entropy  $-\log(Pz(z))$  with constant reconstruction error. The most related works to our approach is image compression with deep learning (Johannes Ballé & Johnston, 2018; Jing Zhou, 2019; Sihan Wen, 2019). In these works, the information entropy of feature map extracted by CNN is minimized. To calculate entropy, Johannes Ballé & Johnston (2018) propose method to estimate probability distribution of latent variable  $Pz(z)$  parametrically, assuming univariate independent (factorized) distribution for each latent variable. We discover that, behind the success of these compression methods, RDO has effect to scale latent variables. Inspired by this, we propose autoencoder which scale latent variables according to definition of distance of data, without limitation of the transformation function. Our scheme is applicable even if GMM is used to estimate  $Pz(z)$ , which is suitable for clustering and anomaly detection. Furthermore, in the case factorized distribution is used for  $Pz(z)$ , our model behaves as continuous PCA. This property is considered to promote the interpretation of latent variable.

## 3 METHOD AND THEORY

The fundamental mechanism of rate-distortion optimization guided auto-encoder is minimizing cost function which consists of (i) reconstruction error between input data and decoder output with noise

to latent variable and (ii) entropy of latent variable. By doing so, the model automatically finds an appropriate scale of latent space. The specific method is described below.

First, let  $\mathbf{x}$  be  $M$ -dimensional domain data ( $\mathbf{x} \in \mathbf{R}^M$ ) and  $P_{\mathbf{x}}(\mathbf{x})$  the probability of  $\mathbf{x}$ . Then  $\mathbf{x}$  is converted to  $N$ -dimensional latent space  $\mathbf{z} \in \mathbf{R}^N$  by encoder. Let  $f_{\theta}(\mathbf{x})$ ,  $g_{\phi}(\mathbf{z})$ , and  $P_{\mathbf{z}\psi}(\mathbf{z})$  be parametric encoder, decoder, and probability distribution function of latent space with parameters  $\theta$ ,  $\phi$ , and  $\psi$ . We assume that each function's parameter is rich enough to fit ideally. Then latent variable  $\mathbf{z}$  and decoded data  $\hat{\mathbf{x}}$  are generated as bellow:

$$\mathbf{z} = f_{\theta}(\mathbf{x}) \quad \hat{\mathbf{x}} = g_{\phi}(\mathbf{z}) \quad (1)$$

Let  $\epsilon \in \mathbf{R}^N$  be noise with each dimension being independent (non-correlated among different dimensions) with an average of 0 and a variance of  $\sigma^2$  as follows:

$$\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_N), E[\epsilon_i] = 0, E[\epsilon_i \cdot \epsilon_j] = \delta_{ij} \cdot \sigma^2 \quad (2)$$

Then, given the sum of latent variable  $\mathbf{z}$  and noise  $\epsilon$ , the decoder output  $\check{\mathbf{x}}$  is obtained as follows:

$$\check{\mathbf{x}} = g_{\phi}(\mathbf{z} + \epsilon) \quad (3)$$

Here, the cost function is defined by Eq. (4)

$$L = -\log(P_{\mathbf{z}\psi}(\mathbf{z})) + \lambda_1 \cdot h(D(\mathbf{x}, \hat{\mathbf{x}})) + \lambda_2 \cdot D(\hat{\mathbf{x}}, \check{\mathbf{x}}) \quad (4)$$

In this equation, the first term is the estimated entropy of the latent distribution.  $D(\mathbf{x}_1, \mathbf{x}_2)$  in the second and the third term is a distance function between  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . It is assumed that this distance function can be approximated by the following quadratic form in the neighborhood of  $\mathbf{x}$ , where  $\delta\mathbf{x}$  is arbitrary micro variation of  $\mathbf{x}$ ,  $\mathbf{A}(\mathbf{x})$  is  $M \times M$  Hermite matrix depending on  $\mathbf{x}$ , and  $\mathbf{L}(\mathbf{x})$  is Cholesky decomposition of  $\mathbf{A}(\mathbf{x})$ .

$$D(\mathbf{x}, \mathbf{x} + \delta\mathbf{x}) \simeq {}^t \delta\mathbf{x} \cdot \mathbf{A}(\mathbf{x}) \cdot \delta\mathbf{x} = \|\mathbf{L}(\mathbf{x}) \cdot \delta\mathbf{x}\|^2 \quad (5)$$

For instance, when  $D(\cdot, \cdot)$  is square error as in Eq. (6),  $\mathbf{A}(\mathbf{x})$  and  $\mathbf{L}(\mathbf{x})$  are Identity matrices.

$$D(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|^2 \quad (6)$$

In the case of SSIM (Zhou Wang (2001)) metric which is close to subjective image quality, a cost  $(1 - SSIM)$  can be also approximated in this quadratic form. This is explained in Appendix C. Let  $h(\cdot)$  in the second term of Eq. (5) be a scale function such as the identity  $h(d) = d$ , the logarithm  $h(d) = \log(d)$ , etc.. The effect of  $h(d)$  is discussed in Appendix D. Then, averaging Eq. (4) according to distributions,  $\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x})$  and  $\epsilon \sim P(\epsilon)$ . By deriving parameters that minimize this value, the encoder, decoder, and probability distribution of the latent space are trained as Eq. (7).

$$\theta, \phi, \psi = \arg \min_{\theta, \phi, \psi} (E_{\mathbf{x} \sim P_{\mathbf{x}}(\mathbf{x}), \epsilon \sim P(\epsilon)} [L]) \quad (7)$$

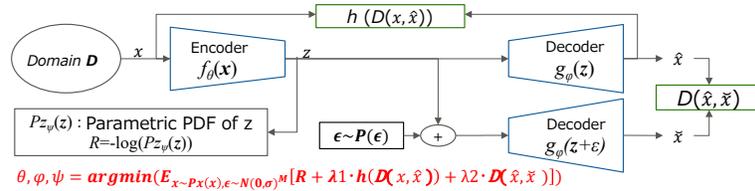


Figure 1: Architecture of RaDOGAGA

Here, we introduce  $\mathbf{x}_D$  as rescaled space of  $\mathbf{x}$  according to the distance function  $D(\cdot, \cdot)$ :

$$d\mathbf{x}_D = \mathbf{L}(\mathbf{x}) \cdot d\mathbf{x} \quad (8)$$

When  $D(\cdot, \cdot)$  is square error, the two spaces of  $\mathbf{x}_D$  and  $\mathbf{x}$  are equivalent. It is turned out that each column of the Jacobian matrix of latent space  $\mathbf{z}$  and rescaled space  $\mathbf{x}_D$  are constant multiple of orthonormal basis regardless of the value of  $\mathbf{z}$  and  $\mathbf{x}$  after training based on Eqs. (4) and (7). Here,  $\delta_{ij}$  denotes Kronecker delta.

$${}^t \left( \mathbf{L}(\mathbf{x}) \cdot \frac{\partial \mathbf{x}}{\partial z_i} \right) \left( \mathbf{L}(\mathbf{x}) \cdot \frac{\partial \mathbf{x}}{\partial z_j} \right) = {}^t \left( \frac{\partial \mathbf{x}_D}{\partial z_i} \right) \left( \frac{\partial \mathbf{x}_D}{\partial z_j} \right) = \delta_{ij} \cdot \frac{1}{2\lambda_2 \sigma^2} \quad (9)$$

A more detailed proof is described in Appendix A. Then, by calculating distance function  $D(\cdot, \cdot)$  for  $\mathbf{x}$  and  $\mathbf{x} + d\mathbf{x}$ , following relationship is derived by Eqs. (5), (8) and (9).

$$D(\mathbf{x}, \mathbf{x} + d\mathbf{x}) = \|\mathbf{L}(\mathbf{x}) \cdot d\mathbf{x}\|^2 = \|d\mathbf{x}_D\|^2 = \left\| \sum_{i=1}^N \frac{\partial \mathbf{x}_D}{\partial z_i} \cdot dz_i \right\|^2 = \frac{1}{2\lambda_2 \sigma^2} \cdot \|dz\|^2 \quad (10)$$

This means that, the latent space is scaled so that the amount of change in the real space in response to the minute change of  $\mathbf{z}$  is constant independent of  $\mathbf{z}$  value.

$$D(\mathbf{x}, \mathbf{x} + d\mathbf{x}) / \|dz\|^2 = \|d\mathbf{x}_D\|^2 / \|dz\|^2 = \text{const.} \quad (11)$$

Next, Jacobian between  $\mathbf{x}_D$  and  $\mathbf{z}$  is examined. First, when  $M = N$ , the Jacobian matrix  $d\mathbf{x}_D/d\mathbf{z}$  is a square matrix, and each column is the same as  $1/(\sqrt{2\lambda_2} \cdot \sigma)$  times orthonormal basis. For this reason, the Jacobian is a constant regardless of the value of  $\mathbf{z}$  as shown below:

$$\left| \frac{d\mathbf{x}_D}{d\mathbf{z}} \right| = \left( \frac{1}{2\lambda_2 \cdot \sigma^2} \right)^{(N/2)} \quad (12)$$

In this case, the probability distribution of  $\mathbf{x}_D$  and  $\mathbf{z}$  is proportional because of the constant Jacobian. For the case of  $M > N$ , we assume the situation where most energy is efficiently and effectively mapped to N-dimensional latent space. Then the product of the singular values of SVD for a Jacobi matrix can be regarded as a pseudo Jacobian between the real space and the latent space. Since all of the N singular values are  $1/(\sqrt{2\lambda_2} \cdot \sigma)$ , the pseudo Jacobian is also a constant. As a result, the following equation holds where  $J$  is Jacobian or pseudo Jacobian.

$$Pz(\mathbf{z}) \simeq J \cdot P\mathbf{x}_D(\mathbf{x}_D) \propto P\mathbf{x}_D(\mathbf{x}_D) \quad (13)$$

Let  $\hat{P}\mathbf{x}_D(\mathbf{x}_D)$  be estimated probability of  $\mathbf{x}_D$ . Because the Jacobian  $J$  is constant,  $\hat{P}\mathbf{x}_D(\mathbf{x}_D)$  can be approximated by  $J^{-1} \cdot Pz_\psi(\mathbf{z})$ . Accordingly, the average of the first term in Equation (4) by  $\mathbf{x} \sim P\mathbf{x}(\mathbf{x})$  can be transformed as follows.

$$E_{\mathbf{z} \sim Pz(\mathbf{z})}[-\log(Pz_\psi(\mathbf{z}))] \simeq -E_{\mathbf{x}_D \sim P\mathbf{x}_D(\mathbf{x}_D)}[\log(\hat{P}\mathbf{x}_D(\mathbf{x}_D))] - \log(J) \quad (14)$$

Here, minimization of Equation (14) is equivalent to the log-likelihood maximization of  $\hat{P}\mathbf{x}_D(\mathbf{x}_D)$ . As a result, the PDF of  $\mathbf{x}_D$  can be estimated without maximizing ELBO.

Regarding a parametric probability distribution, if a model such as GMM is taken, it is considered that a multimodal probability distribution can be obtained flexibly. Besides, when the following factorized probability model is used, the model shows a "continuous PCA" feature where the energy is concentrated in several principal latent variables.

$$Pz_\psi(\mathbf{z}) = \prod_{i=1}^N Pz_{i\psi}(Z_i) \quad (15)$$

The derivation of "continuous PCA" feature is explained in Appendix B.

## 4 EXPERIMENT

### 4.1 PROBABILITY DENSITY ESTIMATION WITH TOY DATA

In this section, we describe our experiment using toy data to demonstrate whether the probability density of the input data  $P\mathbf{x}(\mathbf{x})$  and that of estimated in the latent space  $Pz(\mathbf{z})$  are proportional to each other as in theory. First, we sample data  $\mathbf{s} = (s_1, s_2 \dots s_{10,000})$  from three-dimensional Gaussian distribution consists of three-mixture-components with mixture weight  $\boldsymbol{\pi} = (1/3, 1/3, 1/3)$ , mean  $\boldsymbol{\mu}_k = (\mu_{k1}, \mu_{k2}, \mu_{k3})$ , and covariance  $\boldsymbol{\Sigma}_k = \text{diag}(\sigma_{k1}, \sigma_{k2}, \sigma_{k3})$ .  $k$  is the index for mixture component. Then, we scatter  $\mathbf{s}$  with uniform random noise  $\mathbf{u} \in R^{3 \times 16}$ ,  $u_{dm} \sim U_d(-\frac{1}{2}, \frac{1}{2})$ , where  $d$  and  $m$  are index for dimension of sampled data and scattered data. The  $U_{ds}$  are uncorrelated with each other. We produce 16-dimensional input data  $\mathbf{x}$  with a sample number of 10,000 in the end.

$$\mathbf{x} = \sum_{d=1}^3 \mathbf{u}_d s_d \quad (16)$$

The appearance probability of the input data  $P\mathbf{x}(\mathbf{x})$  is equivalent to the generation probability of  $\mathbf{s}$ .

## 4.1.1 CONFIGURATION

In the experiment, we estimate the  $P_{z_\psi}(z)$  using GMM with parameter  $\psi$  as in DAGMM (Zong et al., 2018). Instead of EM algorithm, GMM parameters are learned using Estimation Network (EN), which consists of multi-layer neural network. When the GMM has  $K$  mixture-components and  $I$  is the size of batch samples, EN outputs the mixture-components membership prediction as  $K$ -dimensional vector  $\hat{\gamma}$  as follows:

$$\mathbf{p} = EN(\mathbf{z}; \psi), \hat{\gamma} = \text{softmax}(\mathbf{p}) \quad (17)$$

$K$ -th mixture weight  $\hat{\phi}_k$ , mean  $\hat{\mu}_k$ , covariance  $\hat{\Sigma}_k$ , and entropy  $R$  of  $\mathbf{z}$  are further calculated by Eqs. (18) and (19).

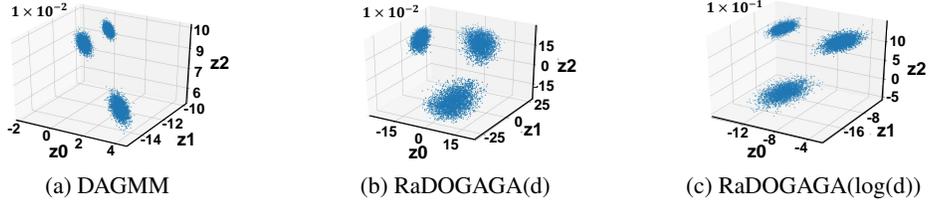
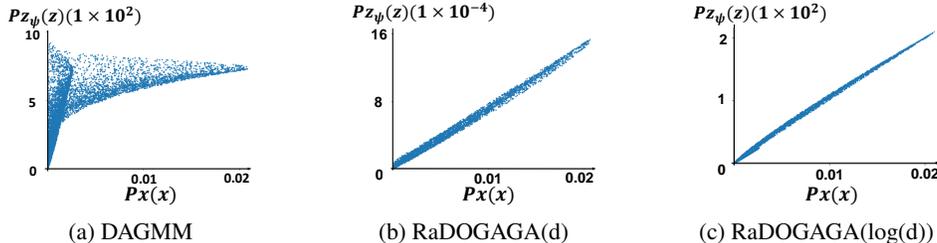
$$\hat{\pi}_k = \sum_{i=1}^I \hat{\gamma}_{ik} / I, \quad \hat{\mu}_k = \sum_{i=1}^I \hat{\gamma}_{ik} \mathbf{z}_i / \sum_{i=1}^I \hat{\gamma}_{ik}, \quad \hat{\Sigma}_k = \sum_{i=1}^I \hat{\gamma}_{ik} (\mathbf{z}_i - \hat{\mu}_k)(\mathbf{z}_i - \hat{\mu}_k)^T / \sum_{i=1}^I \hat{\gamma}_{ik} \quad (18)$$

$$R = -\log \left( \sum_{k=1}^K \hat{\pi}_k / \sqrt{|2\pi\hat{\Sigma}_k|} \cdot \exp \left( -\frac{1}{2} (\mathbf{z} - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1} (\mathbf{z} - \hat{\mu}_k) \right) \right) \quad (19)$$

Overall network is trained by Eqs. (4) and (7). In this experiment, we set  $D(x_1, x_2)$  as square error  $\|x_1 - x_2\|^2$ , and test two types of  $h()$ ,  $h(d) = d$  and  $h(d) = \log(d)$ . We denote models trained with these  $h()$  as RaDOGAGA(d) and RaDOGAGA(log(d)). As a comparison method, DAGMM is used. DAGMM also consists of encoder, decoder, and EN. In DAGMM, to avoid falling into the trivial solution that entropy is minimized when the diagonal component of the covariance matrix is 0, the inverse of the diagonal component is added to the cost function as Eq. (20):

$$L = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \lambda_1 \cdot (-\log(P_{z_\psi}(\mathbf{z}))) + \lambda_2 P(\hat{\Sigma}), \quad P(\hat{\Sigma}) = \sum_{k=1}^K \sum_{j=1}^d \hat{\Sigma}_{kjj}^{-1} \quad (20)$$

For both RaDOGAGA and DAGMM, the autoencoder part is constructed with fully connected (FC) layers with sizes of 64, 32, 16, 3, 16, 32, and 64. For all FC layers except for the last of the encoder and the decoder, we attach  $\tanh$  as the activation function. The EN part is also constructed with FC layer with a size of 10, 3. For the first layer, we attach the  $\tanh$  as activation function and dropout (ratio=0.5). For the last one, softmax is attached.  $(\lambda_1, \lambda_2)$  is set as  $(1 \times 10^{-4}, 1 \times 10^{-9})$ ,  $(1 \times 10^6, 1 \times 10^6)$  and  $(1 \times 10^3, 1 \times 10^4)$  for DAGMM, RaDOGAGA(d) and RaDOGAGA(log(d)) respectively. Optimization is done by Adam optimizer (Kingma & Ba, 2014) with learning rate  $1 \times 10^{-4}$  for all model. We set  $\sigma^2$  as  $1/12$ .

Figure 2: Plot of latent variable  $z$ Figure 3: plot of  $Px(x)$  and  $P_{z_\psi}(z)$

#### 4.1.2 RESULT

Figure 2 displays the distribution of latent variable  $z$ , and Figure 3 displays a plot of  $Px(x)$ (x-axis) against  $Pz_\psi(z)$ (y-axis). Even though both methods can capture that  $s$  is generated from three mixture-component, there is a difference in PDF. In our method,  $Px(x)$  and  $Pz_\psi(z)$  are approximately proportional to each other as in theory while we cannot observe such proportionality in DAGMM. To compare RaDOGAGA(d) and RaDOGAGA(log(d)), we normalized  $Pz_\psi(z)$  with the range from 0 to 1, then calculated residual of linear regression putting  $Pz_\psi(z)$  as objective variable and  $Px(x)$  as explanatory variable. The residual of RaDOGAGA(log(d)), i.e., 0.0027, is smaller than that of RaDOGAGA(d), i.e., 0.0100. Actually, when  $Pz_\psi(z)$  is sufficiently fitted,  $h(d) = \log(d)$  makes  $Px(x)$  and  $Pz_\psi(z)$  be proportional more rigidly. On the other hand,  $h(d) = d$  makes the scale of latent space slightly bent in order to minimize entropy function, allowing relaxed fitting of  $Pz_\psi(z)$ . More detail is described in Appendix D. In the next section, we see how this trait performs on the real task.

### 4.2 ANOMALY DETECTION TASK USING REAL DATA

In this section, we examine whether the clear relationship between  $Px(x)$  and  $Pz(z)$  is useful in the anomaly detection task using real data. We use four public datasets<sup>1</sup>, KDDCUP99, Thyroid, Arrhythmia, and KDDCUP-Rev. The (instance number, dimension, anomaly ratio(%)) of each dataset is (494021, 121, 20), (3772, 6, 2.5), (452, 274, 15), and (121597, 121, 20) respectively. Detail of datasets is described in Appendix E.

#### 4.2.1 EXPERIMENTAL SETUP

We follow the setting in Zong et al. (2018). Randomly extracted 50% of the data is assigned to training and the rest to testing. During training, only normal data is used. During the test, the  $R$  for each sample is calculated as the anomaly score, and if the anomaly score is higher than a threshold, it is detected as an anomaly. The threshold is determined by the ratio of anomaly data in each data set. For example, in KDDCup99, data with  $R$  in the top 20 % is detected as an anomaly. As metrics, precision, recall, and F1 score are calculated. We run 20 times for each dataset split by 20 different random seeds.

#### 4.2.2 BASELINE MODEL

Same as in the previous section, DAGMM is taken as the baseline method. We also compare with the scores reported in previous works in which same experiments were conducted (Zenati et al., 2018; Song & Ou, 2018; Liao et al., 2018).

#### 4.2.3 CONFIGURATION

As in Zong et al. (2018), in addition to the output from the encoder,  $\frac{\|x-x'\|_2}{\|x\|_2}$  and  $\frac{x \cdot x'}{\|x\|_2 \|x'\|_2}$  are concatenated to  $z$ . It is sent to EN. Note that  $z$  is sent to the decoder before concatenation. Other configuration except for hyper parameter is same as in the previous experiment. Hyper parameter for each dataset is described in Appendix E.

#### 4.2.4 RESULTS

Table 1 reports the average scores and standard deviations (in brackets). Comparing DAGMM and RaDOGAGA, RaDOGAGA has a better performance regardless of types of  $h()$ . Simply introducing the RDO mechanism into the autoencoder has a valid efficacy for anomaly detection. Moreover, our approach achieves state-of-the-art performance compared to other previous works in which same datasets is used. Clear relationship between  $Px(x)$  and  $Pz(z)$  by our model is considered to be effective in the task of anomaly detection where the estimating probability distribution is important. In RaDOGAGA, when we compare result of RaDOGAGA(d) and RaDOGAGA(log(d)), either of one is not always superior. As described in section 4.1 and Appendix D,  $h()$  can be an option depending on fitting flexibility of  $Pz(z)$ .

<sup>1</sup>Dataset can be download from (<https://kdd.ics.uci.edu/>) and (<http://odds.cs.stonybrook.edu>)

Table 1: Average and standard deviations(in brackets) of Precision, Recall and F1

Dataset	Methods	Precision	Recall	F1
KDDCup	ALAD*	0.9427(0.0018)	0.9577(0.0018)	0.9501(0.0018)
	INRF*	0.9452(0.0105)	0.9600(0.0113)	0.9525(0.0108)
	GMVAE*	0.952	0.9141	0.9326
	DAGMM*	0.9297	0.9442	0.9369
	DAGMM+ <sup>†</sup>	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))	<b>0.9563(0.0042)</b>	<b>0.9714(0.0042)</b>	<b>0.9638(0.0042)</b>
Thyroid	GMVAE*	<b>0.7105</b>	0.5745	0.6353
	DAGMM*	0.4766	0.4834	0.4782
	DAGMM+ <sup>†</sup>	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)	<b>0.6848(0.0597)</b>	<b>0.6702(0.0585)</b>
Arrhythmia	ALAD*	0.5000(0.0208)	0.5313(0.0221)	0.5152(0.0214)
	GMVAE*	0.4375	0.4242	0.4308
	DAGMM*	0.4909	0.5078	0.4983
	DAGMM+ <sup>†</sup>	0.4985(0.0389)	0.5136(0.0401)	0.5060(0.0395)
	RaDOGAGA(d)	<b>0.5353(0.0461)</b>	<b>0.5515(0.0475)</b>	<b>0.5433(0.0468)</b>
	RaDOGAGA(log(d))	0.5294(0.0405)	0.5455(0.0418)	0.5373(0.0411)
KDDCup-rev	DAGMM*	0.937	0.939	0.938
	DAGMM+ <sup>†</sup>	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))	<b>0.9864(0.0009)</b>	<b>0.9865(0.0009)</b>	<b>0.9865(0.0009)</b>

\*Score of ALAD, INRF, GMVAE and DAGMM is cited from Zenati et al. (2018), Song & Ou (2018), Liao et al. (2018) and Zong et al. (2018) respectively.

<sup>†</sup>DAGMM+ is our implementation. Note that we also test same configuration as in Zong et al. (2018) and achieve similar score as reported (shown in Appendix E).

### 4.3 QUANTIFYING THE IMPACT OF LATENT VARIABLES ON DISTANCE AND BEHAVIOR AS PCA

In this section, we verify that the impact of each latent variable on the distance function can be quantified. As described in section 3, when  $z$  is displaced by a small interval  $\delta$ , the error ratio between two decoder output  $\frac{D(g(z),g(z+\delta))}{\delta^2}$  is constant regardless of the dimension of  $z$ . We verify if this characteristic is observable in a model trained with real data. Once the model is trained, we encode the image  $i$  and obtain  $z_i$ . Then,  $\delta$  is added to  $n$ -th dimension of  $z_i$  and  $\frac{D(g(z_i),g(z_i+\delta))}{\delta^2}$  is calculated for each sample. Finally, the average across all samples is measured. This operation is conducted to each dimension of  $z$  independently. In rest of this section,  $D'(z_n)$  denote  $\frac{D(g(z),g(z+\delta))}{\delta^2}$  for  $n$ -th dimension. We also observe the distribution of  $z$  and how the image looks different in response to the variation of  $z_n$ .

To train model, we use CelebA dataset (Liu et al., 2015), which consists of 202,599 celebrity images. The images are center-cropped so that the image size is 64 x 64.

#### 4.3.1 CONFIGURATION

In this experiment, factorized distributions (Johannes Ballé & Johnston, 2018) is used to estimate  $Pz_\psi(z)$ . For comparison, we evaluate beta-VAE. Both models are constructed with same depth of CNN and FC layers, with 256-dimensional  $z$ . Detail of networks and hyper parameter is written in Appendix F. For RaDOGAGA, we set  $D(x_1, x_2)$  as  $1 - SSIM(x_1, x_2)$  and  $h$  as  $h(d) = d$ . For beta-VAE, reconstruction loss is also  $1 - SSIM(x_1, x_2)$ .  $SSIM(x_1, x_2)$  is defined by Eq. (60).

#### 4.3.2 RESULT

Both models are trained so that the SSIM between input and reconstructed image is around 0.93. Figure 4a and 4b are the variance of latent variables  $z$  arranged in descending order. The red line is the cumulative relative ratio of the variance. In Fig. 4b, variance is concentrated in a specific

dimension. On the other hand, in Fig. 4a, VAE is trained so that each latent variable is fitted to a Gaussian distribution with mean 0 and variance 1, there is no significant difference in the variance of each latent variable. Figures 4c and 4d respectively depicts the average of  $D'(z_n)$  of each of the top nine dimensions with the largest variance of  $z$ . In VAE, the  $D'(z_n)$  varies drastically depending on the dimension  $n$ , while in RaDOGAGA, it is approximately constant regardless of  $n$ . Figure 5 show decoder outputs when individual  $z_n$  is traversed from  $-2\sigma$  to  $2\sigma$ , fixing rest of  $z$  as mean. From the top, each row corresponds to  $z_0, z_1, z_2 \dots$ , and the center column is mean. In Fig. 5b, the image changes visually in any dimension of  $z$ , while in Fig. 5a, depending on  $n$ , there are cases where no significant changes can be seen (such as  $z_4, z_5$ , and so on). This result means that, in RaDOGAGA, the variance of  $z$  directly corresponds to the visual impact and the distance  $D(x_1, x_2)$ , behaving as PCA. Besides, since  $D'(z_n)$  is constant, the variance can be used as a quantifying metric. Although this is unsupervised image reconstruction task, if the task is the semantical problem and distance function is defined so as to solve it, such as a classification, influence of each  $z$  on the semantics is expected to be quantified.

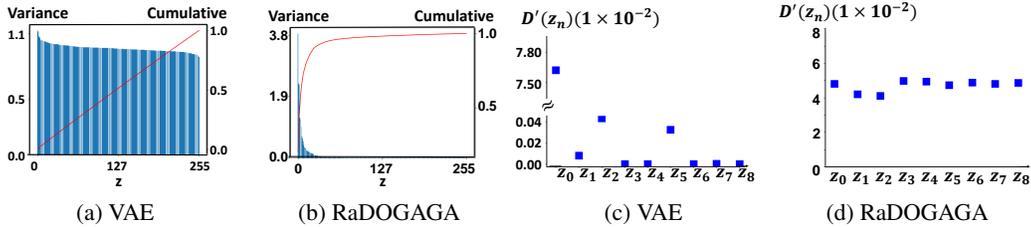


Figure 4: Variance of  $z$  (two on the left) and  $D'(z_n)$  (two on the right)

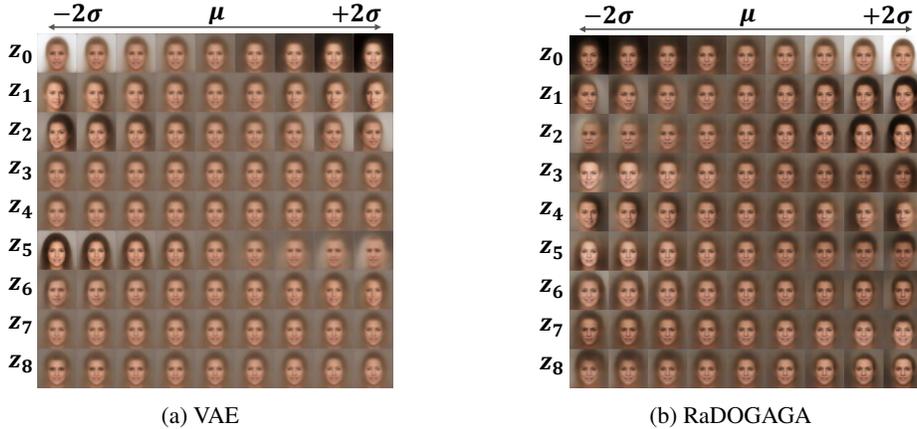


Figure 5: Latent space traversal

## 5 CONCLUSION

In this paper, we proposed RaDOGAGA that learns parametric probability distribution and auto-encoder simultaneously based on rate-distortion optimization. It was proved that the probability distribution of the latent variables obtained by the proposed method is proportional to the probability distribution of the input data theoretically and experimentally. This property is validated in anomaly detection achieving state-of-the-art performance. Moreover, our model has the trait as PCA which likely promotes interpretation of latent variables. For the future work, we will conduct experiments with different types of distance functions that derived from semantical task, such as in categorical classification. Meanwhile, as mentioned in Michael Tschannen & Lucic (2019), considering the usefulness of latent variables in downstream task is another research direction to explore.

## REFERENCES

- Alexander Alemi, Ben Poole, Ian Fischer, Joshua Dillon, Rif A. Saurous, and Kevin Murphy. Fixing a broken elbo. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 159–168. PMLR, 2018.
- Johannes Ballé, Valero Laparra, and Eero P Simoncelli. Density modeling of images using a generalized normalization transformation. *arXiv preprint arXiv:1511.06281*, 2015.
- Tian Qi Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, pp. 2610–2620, 2018.
- Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pp. 2172–2180, 2016.
- Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- A. Pal C. Burgess X. Glorot M. Botvinick S. Mohamed I. Higgins, L. Matthey and A. Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR 2017*, 2017.
- A. Pal C. Burgess X. Glorot M. Botvinick S. Mohamed I. Higgins, L. Matthey and A. Lerchner. H. kim and a. mnih, “disentangling by factorising. In *Proc. of the International Conference on Machine Learning*, pp. 2649–2658, 2018.
- Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou. Variational deep embedding: An unsupervised and generative approach to clustering. *arXiv preprint arXiv:1611.05148*, 2016.
- Akira Nakagawa Kimihiko Kazui Zhiming Tan Jing Zhou, Sihan Wen. Multi-scale and context-adaptive entropy model for image compression. In *In WORKSHOP AND CHALLENGE ON LEARNED IMAGE COMPRESSION*, 2019.
- Saurabh Singh Sung Jin Hwang Johannes Ballé, David Minnen and Nick Johnston. Variational image compression with a scale hyperprior. In *In International Conference on Learning Representations*, 2018.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P. Kingma, Danilo J. Rezende, Shakir Mohamed, and Max Welling. Semi-supervised learning with deep generative models. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’14, pp. 3581–3589, Cambridge, MA, USA, 2014. MIT Press.
- Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10215–10224, 2018.
- Weixian Liao, Yifan Guo, Xuhui Chen, and Pan Li. A unified unsupervised gaussian mixture variational autoencoder for high dimensional outlier detection. In *2018 IEEE International Conference on Big Data (Big Data)*, pp. 1208–1217. IEEE, 2018.
- M. Lichman. Uci machine learning repository. <http://archive.ics.uci.edu/ml/>, 2013.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- Romain Lopez, Jeffrey Regier, Michael I Jordan, and Nir Yosef. Information constraints on auto-encoding variational bayes. In *Advances in Neural Information Processing Systems*, pp. 6114–6125, 2018.

- Olivier Bachem Michael Tschannen and Mario Lucic. Recent advances in autoencoder-based representation learning. In *Third workshop on Bayesian Deep Learning (NeurIPS 2018)*, Montreal, Canada, Dec 2019.
- A.Wiltschko R.P.Adams S.R.Datta M.Johnson, D.K.Duvenaud. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in Neural Information Processing Systems*, 2016.
- Diederik P.Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR 2014*, Banff, Canada, Apr 2014.
- Kamisetty Ramamohan Rao and Pat Yip (eds.). *The Transform and Data Compression Handbook*. CRC Press, Inc., Boca Raton, FL, USA, 2000. ISBN 0849336929.
- Akira Nakagawa Kimihiko Kazui Zhiming Tan Sihan Wen, Jing Zhou. Variational autoencoder based image compression with pyramidal features and context entropy model. In *In WORKSHOP AND CHALLENGE ON LEARNED IMAGE COMPRESSION*, 2019.
- Yunfu Song and Zhijian Ou. Learning neural random fields with inclusive auxiliary generators. *arXiv preprint arXiv:1806.00271*, 2018.
- Aaron van den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, pp. 6306–6315.
- A. Courville Y. Bengio and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35:798–1828, 2013.
- Houssam Zenati, Manon Romain, Chuan-Sheng Foo, Bruno Lecouat, and Vijay Chandrasekhar. Adversarially learned anomaly detection. In *2018 IEEE International Conference on Data Mining (ICDM)*, pp. 727–736. IEEE, 2018.
- Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Balancing learning and inference in variational autoencoders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 5885–5892, 2019.
- H.R. Sheikh E.P. Simoncelli Zhou Wang, A.C. Bovik. *Image quality assessment: from error visibility to structural similarity*. IEEE Trans. on Image Processing, 2001.
- Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International Conference on Learning Representations*, 2018.

## A HOW JACOBI MATRIX BECOME A CONSTANTLY SCALED ORTHONORMAL BASIS

In this appendix section, we prove that all of column vectors in decoder’s Jacobi Matrix have the same norm and are orthogonal to each other. In other words, each column of Jacobi matrix is a constantly scaled orthonormal basis.

Here we assume that data space sample  $\mathbf{x}$  has  $M$  dimension, that is  $\mathbf{x} \in \mathbf{R}^M$ , then encoded to  $N$  dimensional latent space sample  $\mathbf{y} \in \mathbf{R}^N$ . We also assume that an fixed encoder function  $\mathbf{y} = f_{init,\theta}(\mathbf{x})$  and a fixed decoder  $\hat{\mathbf{x}} = g_{init,\phi}(\mathbf{y})$  are given such that  $h(D(\mathbf{x}, \hat{\mathbf{x}}))$  becomes minimal.

$$\mathbf{y} = f_{init,\theta}(\mathbf{x}) \tag{21}$$

$$\hat{\mathbf{x}} = g_{init,\phi}(\mathbf{y}) \tag{22}$$

$$s.t \ h(D(\mathbf{x}, \hat{\mathbf{x}})) \implies \text{minimal}$$

We further assume that fixed parametric PDF of latent variable  $\mathbf{y}$  is also given.

$$P_{\mathbf{y}_{init,\psi}}(\mathbf{y}) \tag{23}$$

Here, it is noted that this function is not needed to be optimal in a sense of  $D_{KL}(P_{\mathbf{y}}(\mathbf{y}) \| P_{\mathbf{y}_{init,\psi}}(\mathbf{y}))$  where  $P_{\mathbf{y}}(\mathbf{y})$  is an actual PDF of  $\mathbf{y}$ .

Under these conditions, we introduce new latent variable  $\mathbf{z}$ , and  $\mathbf{y}$  is transformed from  $\mathbf{z}$  using a following scaling function  $\mathbf{a}(\mathbf{z}) : \mathbf{R}^N \rightarrow \mathbf{R}^N$ .

$$\mathbf{y} = \mathbf{a}(\mathbf{z}) = (a_1(\mathbf{z}), a_2(\mathbf{z}), \dots, a_N(\mathbf{z})) \quad (24)$$

Here, **our goal** is to prove Eq. (9) by examining the condition of this scaling function which minimize the average of Eq. (4) with regard to  $\epsilon \sim P(\epsilon)$ .

Because of the assumption of minimal  $D(x, \hat{x})$ , we ignore the second term of Eq. (4). Next, the PDF of  $\mathbf{z}$  can be derived using a Jacobian of  $\mathbf{a}(\mathbf{z})$ .

$$P_{\mathbf{z}_{\psi,a}}(\mathbf{z}) = \left| \frac{d\mathbf{a}(\mathbf{z})}{d\mathbf{z}} \right| \cdot P_{\mathbf{y}_{init,\psi}}(\mathbf{a}(\mathbf{z})) \quad (25)$$

By applying these conditions and equations to Eq. (4), the cost of scaling function  $\mathbf{a}(\mathbf{z})$  to minimize the average is expressed as follows.

$$L_a = -\log \left( \left| \frac{d\mathbf{a}(\mathbf{z})}{d\mathbf{z}} \right| \cdot P_{\mathbf{y}_{init,\psi}}(\mathbf{a}(\mathbf{z})) \right) + \lambda_2 \cdot D(g_{init,\phi}(\mathbf{a}(\mathbf{z} + \epsilon)), g_{init,\phi}(\mathbf{a}(\mathbf{z}))) \quad (26)$$

Next, the latent variable is fixed as  $\mathbf{y}_0$ , and  $\mathbf{z}_0$  is defined in the next equation.

$$\mathbf{z}_0 = \mathbf{a}^{-1}(\mathbf{y}_0) \quad (27)$$

Then the average of Eq. (26) with regard to  $\epsilon \sim P(\epsilon)$  is expressed as follows.

$$E_{\epsilon \sim P(\epsilon)} [L_a |_{\mathbf{z}=\mathbf{z}_0}] = -\log \left( \left| \frac{d\mathbf{a}(\mathbf{z})}{d\mathbf{z}} \right|_{\mathbf{z}=\mathbf{z}_0} \cdot P_{\mathbf{y}_{init,\psi}}(\mathbf{a}(\mathbf{z}_0)) \right) + \lambda_2 \cdot E_{\epsilon \sim P(\epsilon)} [D(g_{init,\phi}(\mathbf{a}(\mathbf{z}_0 + \epsilon)), g_{init,\phi}(\mathbf{a}(\mathbf{z}_0)))] \quad (28)$$

Then the condition of  $\mathbf{a}(\mathbf{z}_0)$  in the neighborhood of  $\mathbf{y}_0$  is examined when Eq. (28) is minimized. Here, it is noted that the first term of the right side doesn't depend on  $\epsilon$ .

Before examining Eq. (28), some preparation of equations are needed. At first, Jacobi matrix of  $\mathbf{a}(\mathbf{z})$  at  $\mathbf{z} = \mathbf{z}_0$  is defined as  $\mathbf{B}$  using notations of partial differentials  $b_{ij} = \left. \frac{\partial a_i}{\partial z_j} \right|_{\mathbf{z}=\mathbf{z}_0}$ .

$$\mathbf{B} = \left. \frac{\partial \mathbf{y}}{\partial \mathbf{z}} \right|_{\mathbf{z}=\mathbf{z}_0} = \left. \frac{d\mathbf{a}(\mathbf{z})}{d\mathbf{z}} \right|_{\mathbf{z}=\mathbf{z}_0} = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1N} \\ b_{21} & b_{22} & \dots & b_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ b_{N1} & b_{N2} & \dots & b_{NN} \end{pmatrix} \quad \text{where } b_{ij} = \left. \frac{\partial a_i}{\partial z_j} \right|_{\mathbf{z}=\mathbf{z}_0} \quad (29)$$

The vector  $\mathbf{b}_i$  is also defined as follows.

$$\mathbf{b}_i = \left. \frac{\partial \mathbf{y}}{\partial z_j} \right|_{\mathbf{z}=\mathbf{z}_0} = \begin{pmatrix} b_{1i} \\ b_{2i} \\ \vdots \\ b_{Ni} \end{pmatrix} \quad (30)$$

It is clear by definition that  $\mathbf{b}_i$  and  $\mathbf{B}$  have the following relation.

$$\mathbf{B} = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_N) \quad (31)$$

$\Delta_{ij}$  is defined as a cofactor of matrix  $\mathbf{B}$  with regard to the element  $b_{ij}$ , and  $\tilde{\mathbf{b}}_i$  is also defined by the following equation.

$$\tilde{\mathbf{b}}_i = \begin{pmatrix} \Delta_{1i} \\ \Delta_{2i} \\ \vdots \\ \Delta_{Ni} \end{pmatrix} \quad (32)$$

The following equations holds by the definition of cofactor.

$${}^t\mathbf{b}_i \cdot \tilde{\mathbf{b}}_i = \sum_{k=1}^N b_{ki} \cdot \Delta_{ki} = |\mathbf{B}| \quad (33)$$

$$\frac{d|\mathbf{B}|}{db_{ij}} = \Delta_{ij}, \quad \frac{d|\mathbf{B}|}{d\mathbf{b}_i} = \tilde{\mathbf{b}}_i \quad (34)$$

In case of  $i \neq j$ , inner product of  $\mathbf{b}_i$  and  $\tilde{\mathbf{b}}_j$  becomes zero because this value is a determinant of a singular matrix with two equivalent column vectors  $\mathbf{b}_i$ .

$${}^t\mathbf{b}_i \cdot \tilde{\mathbf{b}}_j = \sum_{k=1}^N b_{ki} \cdot \Delta_{kj} = |(\mathbf{b}_1, \dots, \mathbf{b}_i, \dots, \mathbf{b}_i, \dots)| = 0 \quad (35)$$

$\mathbf{G}'_{init}$  is defined as a  $M \times N$  Jacobi matrix of  $g_{init,\phi}(\mathbf{y})$  at  $\mathbf{y} = \mathbf{y}_0$  as follows.

$$\mathbf{G}'_{init} = \left. \frac{dg_{init,\phi}(\mathbf{y})}{d\mathbf{y}} \right|_{\mathbf{y}=\mathbf{y}_0} \quad (36)$$

Using these equations, we proceed to expand of Eq. (28). The first term of (28) in the right side can be expressed as follows, where the second term Eq. (37) in the right side is constant by definition.

$$-\log \left( \left. \frac{d\mathbf{a}(\mathbf{z})}{d\mathbf{z}} \right|_{\mathbf{z}=\mathbf{z}_0} \cdot P\mathbf{y}_{init,\psi}(\mathbf{a}(\mathbf{z}_0)) \right) = -\log(|\mathbf{B}|) - \log(P\mathbf{y}_{init,\psi}(\mathbf{y}_0)) \quad (37)$$

Next step is an expansion of the second term in Eq. (28). First, the following approximation holds.

$$\begin{aligned} g_{init,\phi}(\mathbf{a}(\mathbf{z}_0 + \epsilon)) - g_{init,\phi}(\mathbf{a}(\mathbf{z}_0)) &\simeq \left( \left. \frac{dg_{init,\phi}(\mathbf{y})}{d\mathbf{y}} \right|_{\mathbf{y}=\mathbf{y}_0} \right) \cdot \left( \left. \frac{d\mathbf{a}(\mathbf{z})}{d\mathbf{z}} \right|_{\mathbf{z}=\mathbf{z}_0} \right) \cdot \epsilon \\ &= \mathbf{G}'_{init} \cdot \mathbf{B} \cdot \epsilon \\ &= \mathbf{G}'_{init} \cdot (\epsilon_1 \cdot \mathbf{b}_1 + \epsilon_2 \cdot \mathbf{b}_2 + \dots + \epsilon_N \cdot \mathbf{b}_N) \\ &= \sum_{i=1}^N \epsilon_i \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i \end{aligned} \quad (38)$$

Then, the second term of Eq. (28) can be transformed to the next equation by using Eqs. (38), (2), and (5).

$$\begin{aligned}
E_{\epsilon \sim P(\epsilon)} [D(g_{init,\phi}(\mathbf{a}(\mathbf{z}_0 + \epsilon)), g_{init,\phi}(\mathbf{a}(\mathbf{z}_0)))] &= E_{\epsilon \sim P(\epsilon)} \left[ \left\| \mathbf{L}(\mathbf{x}) \cdot \sum_{i=1}^N \epsilon_i \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i \right\|^2 \right] \\
&= \sum_{i=1}^N E[\epsilon_i^2] \cdot \|\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i\|^2 \\
&\quad + 2 \cdot \sum_{i=1}^N \sum_{j=i+1}^N E[\epsilon_i \cdot \epsilon_j] \cdot \\
&\quad \quad \quad {}^t(\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i) \cdot (\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_j) \\
&= \sigma^2 \cdot \left( \sum_{i=1}^N \|\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i\|^2 \right) \quad (39)
\end{aligned}$$

As a result, Eq. (28) can be rewritten as Eq. (40).

$$\begin{aligned}
E_{\epsilon \sim P(\epsilon)} [L_a | \mathbf{z} = \mathbf{z}_0] &= -\log(|\mathbf{B}|) - \log(P\mathbf{y}_{init,\psi}(\mathbf{y}_0)) \\
&\quad + \lambda_2 \cdot \sigma^2 \cdot \left( \sum_{i=1}^N \|\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i\|^2 \right) \quad (40)
\end{aligned}$$

By examining the minimization process of Eq. (40), the conditions of optimal scaling function  $\mathbf{y} = \mathbf{a}(\mathbf{z})$  in the neighborhood of  $\mathbf{y}_0$  is clarified. Here, the condition of Jacobi matrix  $\mathbf{B}$  is examined instead of  $\mathbf{a}(\mathbf{z})$ . Eq. (40) is differentiated by vector  $\mathbf{b}_i$ , and the result is set to be zero. Then the following equation Eq.(41) is derived.

$$2\lambda_2 \cdot \sigma^2 \cdot ({}^t(\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init}) \cdot (\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init})) \cdot \mathbf{b}_i = \frac{1}{|\mathbf{B}|} \cdot \tilde{\mathbf{b}}_i \quad (41)$$

Afterwards, Eq. (41) is multiplied by  ${}^t\mathbf{b}_j$  from the left, and divided by  $2\lambda_2 \cdot \sigma^2$ . As a result, Eq. (42) is derived by using Eqs. (33) and (35).

$${}^t(\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_j) \cdot (\mathbf{L}(\mathbf{x}) \cdot \mathbf{G}'_{init} \cdot \mathbf{b}_i) = \frac{1}{2\lambda_2 \cdot \sigma^2} \cdot \delta_{ij} \quad (42)$$

Here, we define the following function  $g_{ortho,\psi}(\mathbf{z})$  which is a composite function of  $g_{init,\psi}()$  and  $\mathbf{a}()$ .

$$\hat{\mathbf{x}} = g_{ortho,\psi}(\mathbf{z}) = g_{init,\psi}(\mathbf{a}(\mathbf{z})) \quad (43)$$

Then the next equation holds by definition.

$$\left. \frac{\partial g_{ortho,\psi}(\mathbf{z})}{\partial z_i} \right|_{\mathbf{z}=\mathbf{z}_0} = \left. \frac{\partial \hat{\mathbf{x}}}{\partial z_i} \right|_{\mathbf{z}=\mathbf{z}_0} = \mathbf{G}'_{init} \cdot \mathbf{b}_i \quad (44)$$

It is noted that this equation holds at any value of  $\mathbf{y}_0$  or  $\mathbf{z}_0$ . As a result, the following equation, that is Eq. (9), can be derived.

$${}^t \left( \mathbf{L}(\mathbf{x}) \cdot \frac{\partial \hat{\mathbf{x}}}{\partial z_i} \right) \left( \mathbf{L}(\mathbf{x}) \cdot \frac{\partial \hat{\mathbf{x}}}{\partial z_j} \right) = \delta_{ij} \cdot \frac{1}{2\lambda_2 \sigma^2} \quad (45)$$

If encoder and decoder are trained well and  $\mathbf{x} \simeq \hat{\mathbf{x}}$  holds, we can introduce new rescaled data space  $\mathbf{x}_D$  determined by distance function such as  $d\mathbf{x}_D = \mathbf{L}(\mathbf{x}) \cdot d\mathbf{x}$ , and the next equation holds.

$${}^t \left( \frac{\partial \mathbf{x}_D}{\partial z_i} \right) \left( \frac{\partial \mathbf{x}_D}{\partial z_j} \right) = \delta_{ij} \cdot \frac{1}{2\lambda_2 \sigma^2} \quad (46)$$

In conclusion, all column vectors of Jacobi matrix between  $\mathbf{z}$  and  $\mathbf{x}_D$  has the same L2 norm  $1/\sqrt{2\lambda_2}\sigma$  and all pairs of column vectors are orthogonal. In other words,, when column vectors of Jacobi matrix are multiplied by the constant  $\sqrt{2\lambda_2}\sigma$ , the resulting vectors are orthonormal.

## B EXPLANATION OF "CONTINUOUS PCA" FEATURE

In this section, we explain RaDOGAGA has a continuous PCA feature when factorized probability density model as below is used.

$$Pz_\psi(\mathbf{z}) = \prod_{i=1}^N Pz_{i\psi}(z_i) \quad (47)$$

Here, our definition of "continuous PCA" feature means 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.

Next we explain the reason why these feature is acquired. As explained in appendix A, all column vectors of Jacobi 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 Jacobi 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  $\theta$ , latent variable is described as  $\mathbf{z}_\theta = f_{ortho,\theta}(\mathbf{x})$ , and its PDF is defined as  $Pz_\theta(\mathbf{z}_\theta)$ . PDFs of latent space and data space have the following relation where  $J$  is a Jacobian or pseudo-Jacobian between two spaces with constant value as explained in appendix A.

$$Pz_\theta(\mathbf{z}_\theta) \simeq J \cdot P\mathbf{x}_D(\mathbf{x}_D) \propto P\mathbf{x}_D(\mathbf{x}_D) \quad (48)$$

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

By applying the result of Eqs. (40) and (42), Eq. (4) can be transformed as Eq. (49) where  $\hat{\mathbf{x}} = g_{ortho,\phi}(f_{ortho,\theta}(\mathbf{x}))$ .

$$\begin{aligned} L_{ortho} &= -\log(Pz_\psi(\mathbf{z}_\theta)) + \lambda_1 \cdot \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + N/2 \\ s.t. \quad & \left( \frac{\partial g_{ortho,\phi}(\mathbf{z}_\theta)}{\partial z_{\theta_i}} \right) \cdot \left( \frac{\partial g_{ortho,\phi}(\mathbf{z}_\theta)}{\partial z_{\theta_j}} \right) = \delta_{ij} \cdot \frac{1}{2\lambda\sigma^2} \end{aligned} \quad (49)$$

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

$$L'_{ortho} = -\log(Pz_\psi(\mathbf{z}_\theta)) + \lambda_1 \cdot \|\mathbf{x} - \hat{\mathbf{x}}\|^2 \quad (50)$$

Then the parameters of network and probability can be obtained by the next.

$$\theta, \phi, \psi = \arg \min_{\theta, \phi, \psi} (E_{\mathbf{x} \sim P\mathbf{x}(\mathbf{x})} [L'_{ortho}]) \quad (51)$$

$E_{\mathbf{x} \sim P\mathbf{x}(\mathbf{x})} [L'_{ortho}]$  in Eq. (51) can be transformed as the next.

$$\begin{aligned} E_{\mathbf{x} \sim P\mathbf{x}(\mathbf{x})} [L'_{ortho}] &= \int P\mathbf{x}(\mathbf{x}) \cdot (-\log(Pz_\psi(\mathbf{z}_\theta)) + \lambda_1 \cdot h(D(\mathbf{x}, \hat{\mathbf{x}}))) d\mathbf{x} \\ &= \int \left( Pz_\theta(\mathbf{z}_\theta) \cdot \left| \frac{d\mathbf{x}}{d\mathbf{z}_\theta} \right|^{-1} \right) \cdot (-\log(Pz_\psi(\mathbf{z}_\theta))) \cdot \left| \frac{d\mathbf{x}}{d\mathbf{z}_\theta} \right| dz_\theta \\ &\quad + \lambda_1 \cdot \int P\mathbf{x}_D(\mathbf{x}_D) \cdot |\mathbf{L}(\mathbf{x})|^{-1} \cdot h(D(\mathbf{x}, \hat{\mathbf{x}})) \cdot |\mathbf{L}(\mathbf{x})| d\mathbf{x}_D \\ &= \int Pz_\theta(\mathbf{z}_\theta) \cdot (-\log(Pz_\psi(\mathbf{z}_\theta))) dz_\theta \\ &\quad + \lambda_1 \cdot \int P\mathbf{x}_D(\mathbf{x}_D) \cdot h(D(\mathbf{x}, \hat{\mathbf{x}})) d\mathbf{x}_D \end{aligned} \quad (52)$$

At first, the first term of the third formula in Eq.(52) 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  $\mathbf{z}_{\theta}$ . Then a marginal distribution of  $z_{\theta i}$  can be derived from the next equation.

$$Pz_{\theta i}(z_{\theta i}) = \int Pz_{\theta}(\mathbf{z}_{\theta})d\mathbf{z}_{\theta/i} \quad (53)$$

By using Eqs.(47) and (53), the first term of the third formula in Eq. (52) can be expanded as follows.

$$\begin{aligned} \int Pz_{\theta}(\mathbf{z}_{\theta}) \cdot (-\log(Pz_{\psi}(\mathbf{z}_{\theta}))) d\mathbf{z}_{\theta} &= \int Pz_{\theta}(\mathbf{z}_{\theta}) \cdot \left( -\log \left( \frac{\prod_{i=1}^N Pz_{i\psi}(z_{\theta i})}{\prod_{i=1}^N Pz_{\theta i}(z_{\theta i})} \right) \right) d\mathbf{z}_{\theta} \\ &\quad + \int Pz_{\theta}(\mathbf{z}_{\theta}) \cdot \left( -\log \left( \prod_{i=1}^N Pz_{\theta i}(z_{\theta i}) \right) \right) d\mathbf{z}_{\theta} \\ &= \sum_{i=1}^N \int \left( \int Pz_{\theta}(\mathbf{z}_{\theta})d\mathbf{z}_{\theta/i} \right) \cdot \left( -\log \left( \frac{Pz_{i\psi}(z_{\theta i})}{Pz_{\theta i}(z_{\theta i})} \right) \right) dz_{\theta i} \\ &\quad + \sum_{i=1}^N \int \left( \int Pz_{\theta}(\mathbf{z}_{\theta})d\mathbf{z}_{\theta/i} \right) \cdot (-\log(Pz_{\theta i}(z_{\theta i}))) dz_{\theta i} \\ &= \sum_{i=1}^N D_{KL}(Pz_{\theta i}(z_{\theta i})||Pz_{i\psi}(z_{\theta i})) + \sum_{i=1}^N H(z_{\theta i}) \quad (54) \end{aligned}$$

$H(X)$  is an entropy of variable  $X$ . The first term of the third formula is KL-divergence between marginal probability  $Pz_{\theta i}(z_{\theta i})$  and factorized parametric probability  $Pz_{i\psi}(z_{\theta i})$ . The second term of the third formula can be further transformed using mutual information between latent variables  $I(\mathbf{z}_{\theta})$  and equation (48).

$$\sum_{i=1}^N H(z_{\theta i}) = H(\mathbf{z}_{\theta}) + I(\mathbf{z}_{\theta}) \simeq H(\mathbf{x}_D) - \log(J) + I(\mathbf{z}_{\theta}) \quad (55)$$

The first term of the third formula is the entropy of input data with constant value. The second is also constant. As a result, in order to minimize (54), mutual information  $I(\mathbf{z}_{\theta})$  must be minimized.

At second, the second term of the third formula in Eq. (52) is examined.  $\mathbf{x}_D$  and  $\hat{\mathbf{x}}_D$  denote mapped values to rescaled data space defined by the distance function  $D(\mathbf{x}_1, \mathbf{x}_2)$  as Eq.(8). Because  $\mathbf{x}_D$  and  $\hat{\mathbf{x}}_D$  are close, the following equations holds.

$$\hat{\mathbf{x}}_D \simeq \mathbf{x}_D + \mathbf{L}(\mathbf{x}) \cdot (\hat{\mathbf{x}} - \mathbf{x}) \quad (56)$$

$$D(\mathbf{x}, \hat{\mathbf{x}}) \simeq |\mathbf{L}(\mathbf{x}) \cdot (\hat{\mathbf{x}} - \mathbf{x})|^2 \simeq |\mathbf{x}_D - \hat{\mathbf{x}}_D|^2 \quad (57)$$

By using these expansions, Eq.(52) can be expressed as follows.

$$\begin{aligned} E_{\mathbf{x} \sim P\mathbf{x}(\mathbf{x})}[L'_{ortho}] &\simeq \sum_{i=1}^N D_{KL}(Pz_{\theta i}(z_{\theta i})||Pz_{i\psi}(z_{\theta i})) \\ &\quad + I(\mathbf{z}_{\theta}) + E_{\mathbf{x}_D} \left[ |\mathbf{x}_D - \hat{\mathbf{x}}_D|^2 \right] + \text{Const.} \quad (58) \end{aligned}$$

Here, the rescaled real space  $\mathbf{x}_D$  is divided into a plurality of small subspace partitionings  $\Omega\mathbf{x}_{D1}, \Omega\mathbf{x}_{D2}, \dots$ . Let  $\Omega\mathbf{z}_1, \Omega\mathbf{z}_2, \dots$  be corresponding subspace partitionings in latent space. as the division space of the latent space  $\mathbf{z} \in \mathbf{R}^N$  corresponding to  $\Omega\mathbf{x}_D$ . Then Eq. (58) can be rewritten as follows.

$$\begin{aligned} E_{\mathbf{x} \sim P\mathbf{x}(\mathbf{x})}[L'_{ortho}] &\simeq \sum_{i=1}^N D_{KL}(Pz_{\theta i}(z_{\theta i})||Pz_{i\psi}(z_{\theta i})) \\ &\quad + \sum_k \left( I(\mathbf{z}_{\theta} \in \Omega\mathbf{z}_{\theta k}) + E_{\mathbf{x}_D \in \Omega\mathbf{x}_{Dk}} \left[ |\mathbf{x}_D - \hat{\mathbf{x}}_D|^2 \right] \right) + \text{Const.} \quad (59) \end{aligned}$$

For each subspace partitioning, the transformation from  $\Omega \mathbf{x}_{D_k}$  to  $\Omega \mathbf{z}_{\theta_k}$  can be regarded as constantly scaled orthonormal transformation where orthonormal basis is Jacobi matrix with scale factor  $J^{-1}$ .

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. (59) is minimized, Jacobi matrix from  $\Omega \mathbf{x}_{D_k}$  to  $\Omega \mathbf{z}_{\theta_k}$  for each subspace partitioning should be KLT/PCA. Accordingly, the same feature as PCA will be realized such as the determination of principal components etc.

From these consideration, we conclude that RaDOGAGA has a "continuous PCA" feature.

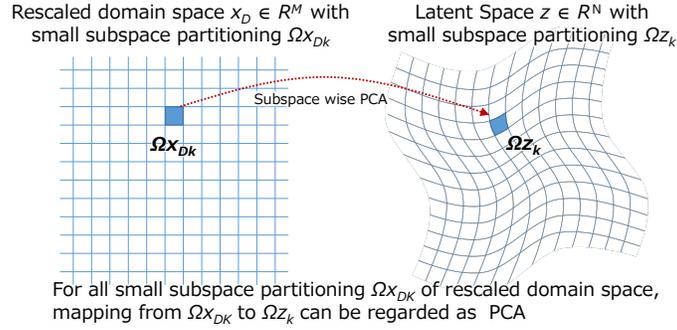


Figure 6: Continuous KLT(PCA) Mapping from input domain to latent space

## C EXPANSION OF SSIM TO A QUADRATIC FORM

Structural similarity (SSIM) (Zhou Wang, 2001) is widely used for picture quality metric which is close to human subjective quality. In this appendix, we show  $(1 - SSIM)$  can be approximated to a quadratic form such as Eq.(5).

Eq. (60) is a SSIM value for a  $N \times N$  window between picture  $\mathbf{x}$  and  $\mathbf{y}$ . In order to calculate SSIM index for a picture, this window is shifted in a whole picture and all of SSIM values are averaged.

$$SSIM_{N \times N}(\mathbf{x}, \mathbf{y}) = \frac{2\mu_x \mu_y}{\mu_x^2 + \mu_y^2} \cdot \frac{2\sigma_{xy}}{\sigma_x^2 + \sigma_y^2} \quad (60)$$

If  $(1 - SSIM_{N \times N}(\mathbf{x}, \mathbf{y}))$  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 approximated as follows

$$SSIM_{N \times N}(\mathbf{x}, \mathbf{x} + \delta \mathbf{x}) = 1 - \frac{\mu_{\delta \mathbf{x}}^2}{2\mu_x^2} - \frac{\sigma_{\delta \mathbf{x}}^2}{2\sigma_x^2} + O\left(\frac{|\delta \mathbf{x}|}{|\mathbf{x}|}\right)^3 \quad (61)$$

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

$$\mu_{\delta \mathbf{x}}^2 = {}^t \delta \mathbf{x} \cdot \mathbf{M} \cdot \delta \mathbf{x}$$

$$\text{where } \mathbf{M} = \frac{1}{N^2} \cdot \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix} \quad (62)$$

$$\sigma_{\delta \mathbf{x}}^2 = {}^t \delta \mathbf{x} \cdot \mathbf{V} \cdot \delta \mathbf{x}$$

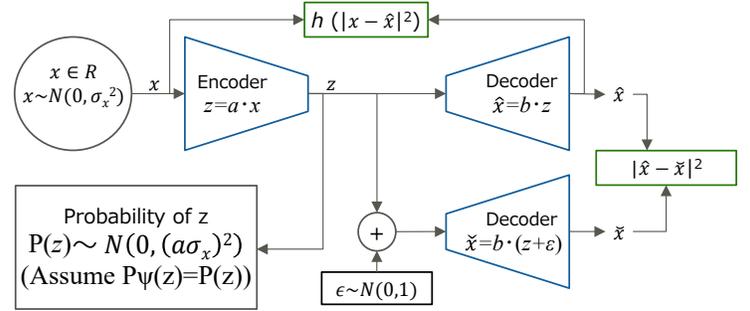
$$\text{where } \mathbf{V} = \frac{1}{N^2} \cdot \begin{pmatrix} N-1 & -1 & \dots & -1 \\ -1 & N-1 & \dots & -1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \dots & N-1 \end{pmatrix} \quad (63)$$

As a result,  $(1 - SSIM_{N \times N}(\mathbf{x}, \mathbf{y}))$  can be expressed in the following quadratic form.

$$1 - SSIM_{N \times N}(\mathbf{x}, \mathbf{x} + \delta \mathbf{x}) \simeq {}^t \delta \mathbf{x} \cdot \left( \frac{1}{2\mu_x^2} \cdot \mathbf{M} + \frac{1}{2\sigma_x^2} \cdot \mathbf{V} \right) \cdot \delta \mathbf{x} \quad (64)$$

## D EFFECT OF $h(x)$

In this appendix section, the effects of two kinds of cost scaling function  $h(d) = d$  and  $h(d) = \log(d)$  are discussed. We evaluated the behaviors of encoder and decoder in a one dimensional model using simple parametric linear encoder and decoder.



$$a, b = \arg \min (E_{x, \epsilon} [-\log(P(z)) + \lambda_1 \cdot h(|x - \hat{x}|^2) + \lambda_2 \cdot |\hat{x} - \tilde{x}|^2])$$

Figure 7: Simple encoder/decoder model to evaluate  $h(d)$

Let  $x$  be a one dimensional data with the normal distribution.

$$x \in R$$

$$x \sim N(0, \sigma_x^2)$$

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

$$z = a \cdot x$$

$$\hat{x} = b \cdot z$$

As a Distance function  $D(x, y)$ , square error is used. The distribution of noise  $\epsilon$  added to latent variable  $z$  is set to  $N(0, 1)$ . Then  $\tilde{x}$  is derived by decoding  $z + \epsilon$ .

$$D(x, y) = |x - y|^2$$

$$\epsilon \sim N(0, 1)$$

$$\tilde{x} = b \cdot (z + \epsilon)$$

For simplicity, we assume parametric PDF  $P\psi(z)$  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.

$$P(z) \sim N(0, (a\sigma_x)^2)$$

$$H(z) = \int -P(z) \cdot \log(P(z)) dz$$

$$= \log(a) + \log(\sigma_x \sqrt{2\pi e})$$

Using these notations, Eqs. (4) and (7) can be expressed as follows.

$$\begin{aligned} Loss &= E_{x \sim N(0, \sigma_x^2), \epsilon \sim N(0, 1)} [-\log P(z) + \lambda_1 \cdot h(|x - \hat{x}|^2) + \lambda_2 \cdot |\hat{x}, \check{x}|^2] \\ &= \log(a) + \log(\sigma_x \sqrt{2\pi e}) + \lambda_1 \cdot E_{x \sim N(0, \sigma_x^2)} [h(|x - \hat{x}|^2)] + \lambda_2 \cdot b^2 \end{aligned} \quad (65)$$

At first, the case of  $h(d) = d$  is examined. By applying  $h(d) = d$ , Eq. (65) can be expanded as follows.

$$Loss = \log(a) + \log(\sigma_x \sqrt{2\pi e}) + \lambda_1 \cdot (a \cdot b - 1)^2 \cdot \sigma_x^2 + \lambda_2 \cdot b^2 \quad (66)$$

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

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

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

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

Here,  $a \cdot b$  is not equal to 1. That is, decoder is not a 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  $P(z) \sim P_\psi(z)$  could be realized at the expense of distance loss or degree of proportional relation  $P(z) \propto P(x)$ .

Next, the case of  $h(d) = \log(d)$  is examined. By applying  $h(d) = \log(d)$  and introducing a minute variable  $\delta$ , Eq. (65) can be expanded as follows.

$$Loss = \log(a) + \log(\sigma_x \sqrt{2\pi e}) + \lambda_1 \cdot \log((a \cdot b - 1)^2 + \delta) + \lambda_2 \cdot b^2 \quad (67)$$

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.

$$\begin{aligned} a \cdot b &= 1 \\ a &= \sqrt{2 \cdot \lambda_2} \\ b &= 1/\sqrt{2 \cdot \lambda_2} \end{aligned} \quad (68)$$

Here,  $a \cdot b$  is equal to 1 and decoder becomes a inverse function of encoder regardless of the variance  $\sigma_x^2$ . In this case, good proportional relation  $P(z) \propto P(x)$  could be realized regardless of the fitting  $P_\psi(z)$  to  $P(z)$ .

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

## E DETAIL OF THE EXPERIMENT IN SECTION 4.2

In this section, we provide further detail of experiment in section 4.2. First, we describe the detail of following four public datasets:

**KDDCUP99 (Lichman (2013))** 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 labeled -normal- and the rest labeled as -attacks-, -normal- instances are used as anomalies, since they are in a minority group.

**Thyroid (Lichman (2013))** 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 (Lichman (2013))** 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 (Lichman (2013))** 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.

Next, hyper parameter for RaDOGAGA is described in table 2. First and second column is number of neuron. For DAGMM, we set same number of neuron in table 2 and  $(\lambda_1, \lambda_2)$  as (0.1, 0.005). Optimization is done by Adam optimizer with learning rate  $1 \times 10^{-4}$  for all dataset.

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	100	10000	100	1000
Arrhythmia	10, 4, 10	10, 2	1000	100	1000	100
KDDCup-rev	60, 30, 8, 30, 60	10, 2	1000	100	100	100

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(referring (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.

## F DETAIL OF THE EXPERIMENT 4.3

In this section, we provide further detail of experiment in section 4.3. For both RaDOGAGA and beta-VAE, we first extract feature with following Convolution Neural Network(CNN).

CNN(9, 9, 2, 64, GDN)-CNN(5, 5, 2, 64, GDN)-CNN(5, 5, 2, 64, GDN)-CNN(5, 5, 2, 64, GDN).

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)

Here, CNN(w, h, s, c, f) is a CNN layer with kernel size (w, h), stride size s, dimension c, and activate function f. GDN(Ballé et al. (2015)) is often used in image compression. Then, we reshape feature map and send to autoencoder as follows.

FC(1024, 8192, softplus)-FC(8192, 256, None)-FC(256, 8192, softplus)-FC(256, 1024, softplus)

FC(i, o, f) is FC layer with input dimension i, output dimension o and activate function f. None means no activate function. Note that, for beta-VAE, since it produces mean and variance, the bottom of the encoder has 2 branches.

$(\lambda_1, \lambda_2)$  is as (1.0, 0.1) is set for RaDOGAGA and  $\beta$  is set as  $1 \times 10^{-4}$  for beta-VAE.

Optimization is done by Adam optimizer with learning rate  $1 \times 10^{-4}$ .