# Latent Feature Decentralization Loss for One-Class Anomaly Detection

**EUNGI HONG**, (Graduate Student Member, IEEE),
**AND YOONSIK CHOE**, (Senior Member, IEEE)
Department of Electrical and Electronic Engineering, Yonsei University, Seoul 03722, South Korea

Corresponding author: Yoonsik Choe (yschoe@yonsei.ac.kr)

**ABSTRACT** Anomaly detection is essential for many real-world applications, such as video surveillance, disease diagnosis, and visual inspection. With the development of neural networks, many neural networks have been used for anomaly detection by learning the distribution of normal data. However, they are vulnerable to distinguishing abnormalities when the normal and abnormal images are not significantly different. To mitigate this problem, we propose a novel loss function for one-class anomaly detection: decentralization loss. The main goal of the proposed method is to cause the latent feature of the encoder to disperse over the manifold space, such that the decoder can generate images similar to those in a normal class for any input. To this end, a decentralization term designed based on the dispersion measure for latent vectors is also added to the existing mean-squared error loss. To design a general solution for various datasets, we restrict the latent space by designing a decentralization loss term-based upper bound of the dispersion measure. As intended, a model trained with the proposed decentralization loss function disperses vectors on the manifold space and generates constant images. Consequently, the reconstruction error increases when the given test image is unknown. Experiments conducted on various datasets verify that the proposed function improves detection performance improved by about 1 % while reducing training time by 48 %, without any structural changes in the conventional autoencoder.

**INDEX TERMS** Anomaly detection, autoencoder, generative network, neural network, representation learning.

## I. INTRODUCTION

Anomaly detection involves the identification of unusual patterns of data not observed during the training phase. It has been continuously studied, owing to its versatile applicability for solving various real-world problems (e.g., video surveillance, disease diagnosis, and visual inspection). An underlying assumption for anomaly detection is that abnormal samples differ from normal samples in both high- and low-dimensional space. Therefore, researchers have focused on mapping techniques to project differences in high-dimensional space that are well-represented in low-dimensional space.

Research on anomaly detection has been conducted using conventional techniques, including principal component

The associate editor coordinating the review of this manuscript and approving it for publication was Hossein Rahmani.

analysis, one-class support vector machines, clustering, and hidden Markov models [2]–[7]. However, with the significant progress made by deep neural networks for many applications in the field of computer vision, many researchers have exploited neural networks as a mapping tool for converting high-dimensional images to low-dimensional latent vectors, because they can non-linearly map high-dimensional data to simple distributions in low-dimensional space. Nevertheless, the number of data may be insufficient to model the statistical characteristics of normal and abnormal data, because it is difficult to obtain samples for abnormalities. Consequently, the training data are usually configured only for normal samples. Therefore, most studies have applied generative networks, such as autoencoders, variational autoencoders (VAEs), adversarial autoencoders, and generative adversarial networks (GANs), in an unsupervised manner [8]–[10]. Most existing methods attempt to train a model to generate realistic

images in various ways, leading the model to learn the distribution of the normal data. This implies that most methods rely on the assumption that abnormal samples are far from normal in the latent vector space. Therefore, models that learn only from normal samples have difficulty reconstructing abnormal data. In other words, they have a high reconstruction error. However, in practice, abnormal samples are not very different from normal samples even in high-dimensional space; they have a similar overall structure with a partial difference. The distance between normal and abnormal data on the latent space is not sufficiently large, leading to similarities in the reconstructed images of abnormal samples and the input images. Consequently, only marginal differences exist in the reconstruction errors between cases where the input is a normal sample and those where it is an abnormal sample.

Motivated by this observation, we approach the problem from a different angle and propose a new loss function for one-class anomaly detection to maximize reconstruction errors for abnormal samples. We use the mean and variance as the central tendency and dispersion measures, respectively, and we confine the space of the latent vector to design a loss function that can find a general solution regardless of the statistical characteristics of the dataset. By restricting the latent space, we obtain the upper limit of the variance of the latent vectors not affected by the characteristics of the dataset. Additionally, the mean vector is continuously updated over iteration steps. Minimizing the proposed loss function containing the decentralization term, we force the encoder to disperse the latent vector for the normal class into broad regions of the manifold space. Simultaneously, the decoder is trained to reconstruct the vectors from the encoder into an image characterized by a normal class, even when fed abnormal samples.

We experimented with proposed algorithm on the MNIST, Fashion MNIST, and MVTEC anomaly detection datasets [11]. Through various experiments, we prove that proposed algorithm has better performance than the existing algorithms; the effectiveness of the proposed algorithm can be confirmed qualitatively through the figures in this paper.

Our contributions can be summarized as follows:

- Instead of focusing on expressing the features of a normal class well by simply reducing the within-class variation, we consider degrading the reconstructed images of abnormal samples. We design a decentralization term with a concept one step beyond that used in previous studies. Through this term, robust anomaly detection is possible, even if the normal and abnormal data are not clearly clustered.
- We design a loss term that is much simpler and easier to implement than the existing method while maintaining state-of-art performance. Furthermore, the designed regularization term can be applied to various datasets by setting the upper bounds of the within-class variance according to the size of the latent space. As a result, more efficient anomaly detection is possible.

- Through various experiments, we show that the proposed algorithm exhibits better performance than the state-of-the-art algorithms, and it generates images having the characteristics of the normal samples for any input by spreading features across latent space. It polarizes the anomaly score for normal and abnormal samples. The proposed algorithm is designed for a unimodal case. However, there are no side effects in the multi-model case. Thus, it can be used for the attention module in multi-class classification, out-of-distribution, and open-set recognition [12], [13].

The remainder of this paper is organized as follows. Section II gives an overview of related work on anomaly detection. Section III elaborates on the proposed method incorporating the decentralization term. Section IV reports and analyzes the experimental results obtained. Finally, the conclusion are given in Section V.

## II. RELATED WORK

Recent developments in neural networks have led to significant progress in supervised learning tasks in computer vision. Various neural network models, including convolutional neural networks, GANs, VAEs, and the adversarial autoencoder have also been used to detect anomalies [15].

The autoencoder first projects the training data onto a low-dimensional space and then inverse-projects them onto the high-dimensional space. The reconstruction error, which is the difference between the input and reconstructed image, is a measure of the difference between normal and abnormal samples. A high reconstruction error means that the input is far from the normal samples. An autoencoder primarily aims at reducing reconstruction errors as an objective function. Most common methods calculate the reconstruction error as mean-squared error(MSE). Bergmann *et al.* considered the structural similarity measure instead of the MSE [16]. The structural similarity measure helps capture the inter-dependencies of adjacent pixels because it considers three different statistical measures: luminance, contrast, and structure. However, in practice, autoencoders tend to yield blurred reconstructions because they regress mostly the conditional mean, rather than the actual multimodal distribution. VAEs mitigate this problem by learning a mapping to a low-dimensional representation, where the actual distribution is modeled. An and Cho showed that the probabilistic characteristics of VAEs aided anomaly detection [17]. They stated that the reconstruction probability was a much more objective and principled anomaly score than the reconstruction error.

Among deep learning methods, GANs have attracted considerable attention owing to their state-of-the-art performance in modeling complex high-dimensional image distributions. Consequently, GANs have been widely used for anomaly detection. Schlegl *et al.* exploited GANs pre-trained for normal samples and detected abnormal samples located far from normal samples in latent space [18]. Zenati *et al.* learned two networks simultaneously to make the whole process more efficient [19]. Akcay *et al.* proposed a network
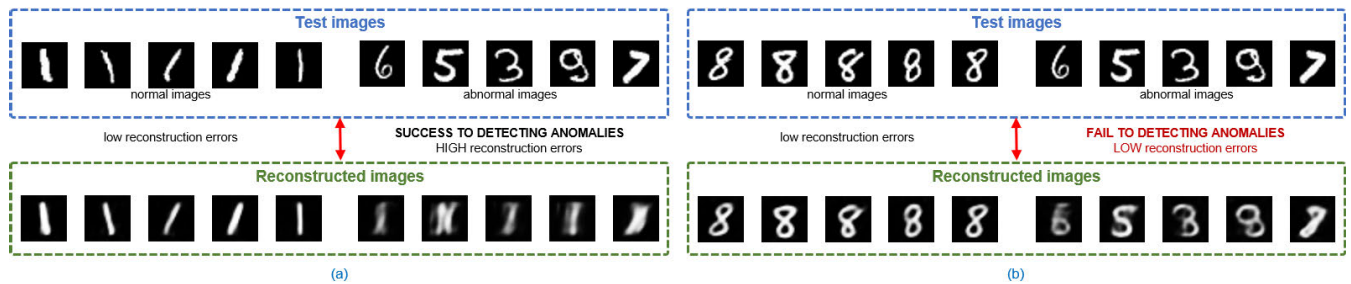
**FIGURE 1.** Illustration of the limitation of the existing method: (a) input and reconstructed images by autoencoder [1] trained for MNIST class "1" as normal, (b) input and reconstructed images by autoencoder [1] trained for MNIST class "8" as normal.
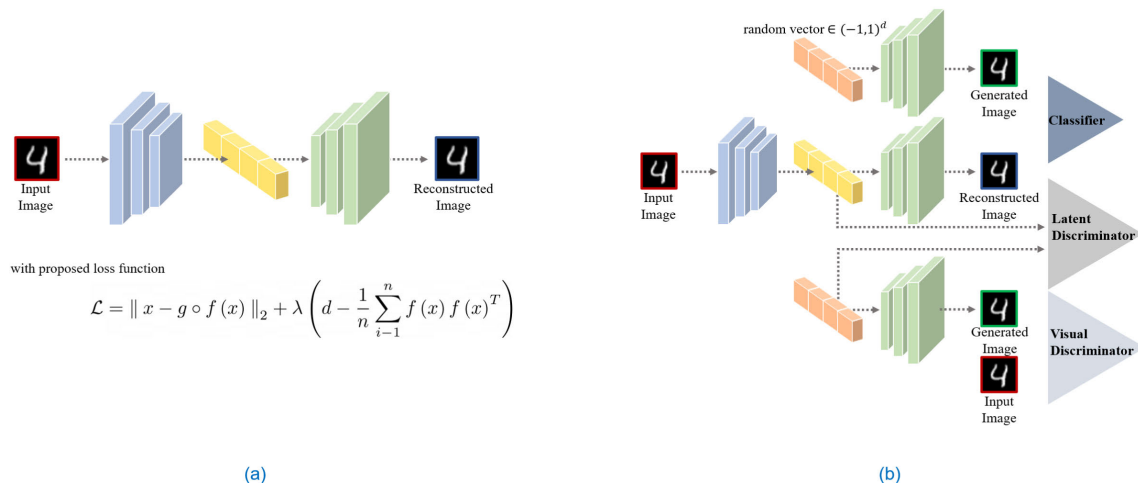


**FIGURE 2.** Illustration of the proposed method and OCGAN [14]: (a) network architecture of proposed method, (b) network architecture of OCGAN.

comprising an encoder–decoder–encoder structure with an adversarial learning scheme to capture the distribution of normal samples [20]. Sabokrou *et al.* trained two modules (i.e., reconstructor and discriminator) via adversarial learning to reconstruct more realistic images Akcay *et al.* exploited adversarial learning over a skip-connected encoder–decoder network architecture [21]. Akcay *et al.* exploited adversarial learning over a skip-connected encoder–decoder network architecture [22]. Skip-connected generator networks capture the details of images well and reconstruct high-quality images drawn from the distribution the model has learned. Ngo *et al.* proposed the Fence-GAN method, which attempts to teach a model the boundary of the normal data distribution [23]. They designed encirclement and dispersion losses to generate data located on the boundary of the normal data distribution instead of overlapping with the data distribution.

All these methods attempted to induce a model to learn a good latent representation that preserves the characteristics of normal samples. However, they can cause a model to reconstruct an image similar to the unknown input. Fig. 1 shows the difference between output images from the model trained on the Modified National Institute of Standards and Technology database (MNIST) class "1" as a normal and output images

from the model trained for class "8" as a normal. When the model is trained for MNIST class "1" as a normal, which is distinctly different from the other classes, reconstructed images for abnormal classes are degraded. However, when learning class "8," which is not clearly distinguished from other classes, the model represents an input-like result image for abnormal samples. Consequently, it is difficult to distinguish abnormalities, owing to a small restoration error. To tackle to this issue, we propose a loss function to spread the feature.

Perera *et al.* raised the issue of previous studies and proposed one-class novelty detection using GAN (OCGAN) [14]. the architecture of which comprises an autoencoder, two discriminators for the latent vector and images, and a classifier. They designed the loss function such that the result of each discriminator for randomly generated vectors from a limited manifold space is always a normal. OCGAN changes abnormal inputs to normal images. As shown in Fig. 2, however, the overall architecture of OCGAN is very complex. Furthermore, adversarial learning-based training is a complicated procedure.

Therefore, we propose a simple and powerful method for anomaly detection via the redesign of loss function for one-class anomaly detection.
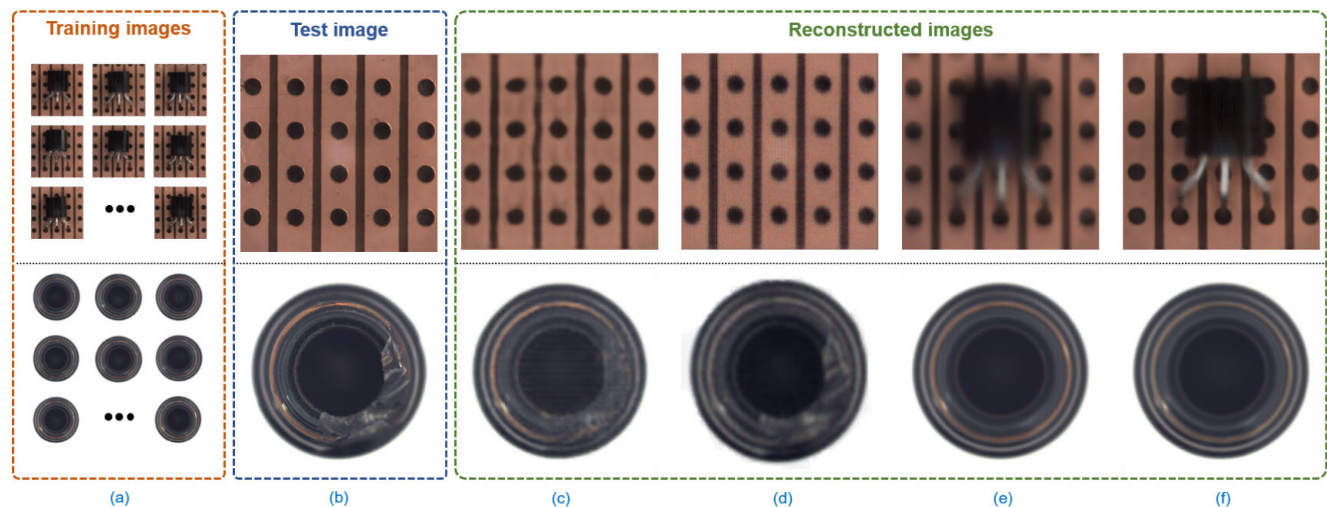
**FIGURE 3.** Illustration of the limitation of existing methods for transistor (top) and bottle (bottom) classes in the MVTec dataset. (a) Normal images used in the training phase, (b) the abnormal image used for testing, (c) the image reconstructed by the autoencoder, (d) the image reconstructed by GANomaly [20], (e) the image reconstructed by OCGAN [14], and (f) the image reconstructed by the proposed method.

## III. PROPOSED METHOD

### A. PROBLEM DEFINITION

We consider a training dataset $X = \{x_1, \ldots, x_n\}$ comprising $n$ normal samples from one-class and an autoencoder model $A$ with an encoder $f$ and a decoder $g$. Model $A$ learns the distribution of $X$ by minimizing $\mathcal{L}_{Recon}$, the difference between the model's input image $x$ and the output image:

$$\mathcal{L}_{Recon} = \| x - g \circ f(x) \|_2 . \quad (1)$$

During the inference phase, for a given test image, an anomaly score $\epsilon$ can be calculated as follows:

$$\epsilon = \| x - g \circ f(x) \|_2 . \quad (2)$$

A high anomaly score indicates that the given data are anomalies. Whereas it is essential to learn the characteristics of normality data by having the loss functions that train a model reconstruct similar images to the input, maximizing the differences between the reconstruction error of normal and abnormal inputs should also be considered. In terms of information theory, training an autoencoder to minimize (1) is equivalent to maximizing the lower bound of mutual information between a high-dimensional image x and a low-dimensional representation $f(x)$ as follows [1]:

$$\arg \max_{\theta} \mathbb{E}[\log p(x|y = f(x)); \theta]. \quad (3)$$

where $\theta$ represents a parameter of the model. It signifies that a network well-trained via (1) effectively reconstructs its inputs from normal samples but is ineffective for reconstructing images from abnormal samples. In real-world applications, however, the between-class variation of the normal and abnormal class in the high-dimensional image space is not sufficiently large, and the anomaly-class image has the same structural characteristics as the normal samples. This leads to there being no difference between abnormal and normal latent

features in the manifold space. Therefore, it is impossible to solve the problem by minimizing within-class variation of features for the normal class. In this case, the model will generate an image that is the same as the input image, even for the anomaly class. Therefore, it has a small restoration error, as can be seen in Fig. 3.

To maximize the anomaly score for abnormal samples, we focus on ensuring that the model can generate a normal class image at any time. In other words, it is necessary to learn to emphasize the difference between a learned class and a non-learned class, rather than only focusing on effectively reconstructing the learned class. In the following sub-section, we introduce a new loss function that implements this concept. In the following sub-section, we introduce a new loss function that implements this concept.

### B. PROPOSED DECENTRALIZATION LOSS

We developed a novel loss function to improve the discriminative power of the MSE-based anomaly score. The objective of the loss function is to maximize the anomaly score for abnormal samples while maintaining a small anomaly score for normal samples. To achieve this, we use an approach that induces the network to always generate an image of the trained class. To this end, based on the assumption that the decoder reconstructs images well only for the trained latent feature, we add the regularization term to the loss function (1) to spread vectors for the normal class across the entire manifold space by maximizing the dispersion measure between latent vectors of the normal class and its central tendency, which is the central value of the distribution. We call this term the decentralization term. This allows the latent vector to be located over a broader range of the manifold space. Through this term, a decoder generates constant images with the characteristics of the normal class for any input image.

### 1) CENTRAL TENDENCY AND DISPERSION MEASURE

In one-class anomaly detection, it is not necessary to consider the covariance of other classes. Therefore, the Euclidean distance is a suitable distance measure. Furthermore, derivations of the $l_2$-norm are easily computed. It is also easy to use gradient-based learning methods. Thus, the loss function should be designed based on the $l_2$-space. To maximize the distance between the latent vectors in the $l_2$-space, we optimize the loss function through the central tendency and dispersion measure of the $l_2$-space. In this paper, $C$ denotes the central tendency of a normal class, and we use the mean vector as a central tendency. For a given dataset $X$, thought of as a vector $f(x) = (f_1(x), \ldots, f_d(x))$, where $d$ is the dimension of the vector, dispersion measure about a central tendency $C$ is the distance from $f(x)$ to the mean vector $C$ in the $p$-norm as follows:

$$\mathbb{D}_p(C) = \|f(x) - C\|_p := \left( \frac{1}{d} \sum_{i-1}^{d} |f_i(x) - C|^p \right)^{1/p}. \quad (4)$$

According to (4), the dispersion measure $\mathbb{D}_2(C)$ becomes the standard deviation and can be replaced by the variance term owing to its proportional property. Therefore, we design a loss function that maximizes the variance.

### 2) ROBUST DECENTRALIZATION LOSS

As mentioned, the objective of the proposed method is to maximize the anomaly score $\epsilon$ in (2) for abnormal samples. The key is to cause the model to reconstruct a representative image with the structural characteristics of the normal samples.

To this end, we designed a loss function, referred to as decentralization loss, that can maximize the variance of the distribution of feature vectors. However, because the value of the latent vector can have an infinite range, it is impossible to apply it directly to the objective function. Even if the reciprocal of the variance is used as a solution to this, there is no general solution. This is because variance varies up to four times or more, depending on the dataset, although there is no different of the result for (1). Therefore, the parameter for adjusting the balance of (1) and the decentralization term is required and must be changed according to the dataset. To determine such a parameter, various factors should be considered, including the within-class and between-class variations of the dataset. However, because it is difficult or even impossible to consider these factors in one-class abnormal detection, we confine the value of the latent vector through the activation function. Then, the decentralization term can be designed based on the upper limit of the variance. The decentralization loss can be expressed as follows:

$$\mathcal{L}_D = d - (C - f(x)) (C - f(x))^T. \quad (5)$$

Here, encoded output $f(x)$ can be represented as $f(x) \in (-1, 1)^d$, where $d$ is the dimension of the latent space. Since we use a tangent hyperbolic function as the activation function, the maximum value of the variance is the same as the

size of the latent vector $d$. The upper limit can be calculated using Theorem 1. and Corollary 1. Assuming that the value of the latent vector has a limited range, the upper limit of the variance of the distribution can be defined by various inequalities—such as Popoviciu's inequality on variance, which is an upper bound on the variance of any bounded probability distribution [24].

*Theorem 1:* Bhatia-Davis inequality [25].

Suppose a distribution has a minimum $m$, a maximum $M$, and an expected value $\mu$. Then, according to the Bhatia–Davis inequality,

$$\frac{1}{d} \sum_{i-1}^{d} |C - f_i(x)|^2 \leq (M - C)(C - m). \quad (6)$$

*Corollary 1:* Popoviciu's inequality on variance.

Let $M$ and $m$ be the upper and lower bounds on the values of any random variable with a particular probability distribution. Then, according to Popoviciu's inequality

$$\frac{1}{d} \sum_{i-1}^{d} |C - f_i(x)|^2 \leq \frac{1}{4}(M - m)^2. \quad (7)$$

Therefore, it is possible to calculate the upper bound for the variance of the given data in the manifold space.

In the next subsection, we design a total loss function containing the proposed decentralization loss.

### 3) TOTAL LOSS FUNCTION

If we minimize the decentralization loss only, the latent features will disperse throughout the space, and the reconstructed images will be degraded. However, if we use only the MSE loss, the resulting latent features would have small within-class variations, and the reconstructed images will be similar to the inputs for anomalies. This means that using one technique or the other is not suitable for anomaly detection as much as using both techniques. Therefore, it is essential to combine these two components of loss, optimize them simultaneously, and balance the two objectives, as confirmed through experiments. Then, the loss function can be expressed as follows:

$$\mathcal{L} = \mathcal{L}_{Recon} + \lambda \cdot \mathcal{L}_D, \quad (8)$$

where the decentralization loss serves as a regularization term and $\lambda$ is a regularization parameter. By limiting the range of latent feature values, we can obtain the upper limit of variance of latent feature vectors. As a result, the decentralization loss has a value in a certain range regardless of the dataset, and it is possible to set a $\lambda$ value applicable to all datasets. We verified this via experimentation, and the optimal $\lambda$ value was determined to be 0.01. The whole loss function can be expressed as follows:

$$\mathcal{L} = \frac{1}{m} \sum_{i-1}^{m} \| x_i - g \circ f(x_i) \|_2$$
$$+ \lambda \cdot \left( d - (C - f(x_i)) (C - f(x_i))^T \right). \quad (9)$$

where $m$ represents size of mini-batch. To effectively maximize the variations, the central tendency $C$ should be updated as the latent features change. Thus, the latent vectors of the entire training set should be considered in every iteration to calculate the central tendency of the normal data, which is inefficient and even impractical. Therefore, the loss function containing the central tendency cannot be used directly before the center loss [26]. To address this problem, instead of updating the centers with respect to the entire training set, we perform the update based on mini-batches. In each iteration, the central tendency is computed by averaging the features of the corresponding classes. The gradients of $\mathcal{L}_D$ with respect to $f(x)$ are computed and the equation of $C$ is updated and computed as follows:

$$\frac{\partial \mathcal{L}_D}{\partial f(x)} = C_j(f(X)) - f(x), \quad (10)$$

$$\Delta C_j(f(X)) = \frac{-\sum_{i=1}^{m}(f(x_i) - C_j(f(X)))}{m+1}. \quad (11)$$

where $j$ represents the iteration number. Additionally, we want to further emphasize one of the strengths of the proposed method: easy implementation. Thus, we devised a simplified version of the loss function. Existing studies have stated that the nonlinearity of neural networks is capable of projecting data into a specific distribution. Thus, we should be able to guide the central points of the class from which we want the model to learn [27]–[32]. By fixing the central point, we can apply Theorem 1, which is stronger than Popoviciu's inequality on variances. By setting the mean vector to zero, the decentralization term can be simply expressed as follows:

$$\mathcal{L}_D = d - \frac{1}{m} \sum_{i-1}^{m} f(x_i)f(x_i)^T. \quad (12)$$

Then, the objective function can be simply expressed as follows:

$$\mathcal{L} = \frac{1}{m} \sum_{i-1}^{m} \| x_i - g \circ f(x_i) \|_2 + \lambda \left( d - f(x_i)f(x_i)^T \right). \quad (13)$$

As a result, we can emphasize the easy implementation, a strength of the proposed algorithm, while maintaining performance. We have demonstrated this through various experiments.

## IV. EXPERIMENTAL EVALUATION
### A. IMPLEMATION DETAILS
#### 1) ENVIRONMENTS
Our algorithm was implemented using PyTorch 1.2.0, and all experiments were conducted on a computer equipped with an Intel i7-9700 processor, 32-GB RAM, and two RTX 2080Ti 11-GB graphical processing units.

#### 2) NETWORK ARCHITECTURE AND HYPER-PARAMETER
##### a: MNIST/FASHION MNIST
We used the same autoencoder architecture as that of OCGAN for MNIST and Fashion MNIST. The autoencoder

was a symmetric network with three $5 \times 5$ convolutions with a stride of two, followed by three transposed convolutions. All convolutions and transposed- convolutions were followed by batch normalization operations and a leaky rectified linear unit (ReLU) having a slope of 0.2. A tanh activation was placed immediately after the last convolution layer to restrict latent-feature values. The initial number of channels was 64. The input and output size were $28 \times 28 \times 1$. Training epochs were 200, and the regularization parameter, $\lambda$, was 0.01.

##### b: MVTec
For the MVTec dataset, the autoencoder was also a symmetric network having eight $4 \times 4$ convolutions and a stride of two, followed by eight transposed convolutions. All convolutions and transposed convolutions were followed by batch normalization operations and a leaky ReLU having a slope of 0.2. The activation function was tanh. The initial number of channels was 64. The input and output size were $256 \times 256 \times 3$. Training epochs were 200, and the regularization parameter, $\lambda$, was 0.01.

### B. METRICS
Performance comparisons were made considering the area under the receiver operating characteristic (AUROC), which is a performance evaluation method that considers the true-and false-positive rates. The performance was also compared in terms of the average and variance of the AUROC obtained from the same five experiments for more accurate performance measurements.

### C. DATASETS
The most widely used datasets (i.e., MNIST and Fashion MNIST) were used for the comparison with other methods. Because these datasets are not designed for anomaly detection, we trained the model for only one class as a normal class, and the performance was evaluated for the entire test dataset. Classes other than the trained class were assumed to be abnormal. To experiment even when the normal and abnormal data had a similar overall structure, we also conducted experiments on the MVTec dataset, which is similar to the actual anomaly detection situation.

### D. EXPERIMENTAL RESULTS
#### 1) AUROC RESULTS
##### a: MNIST
The MNIST dataset, having classes "0" to "9" and a resolution of $28 \times 28$, is the most widely used dataset for one-class anomaly detection. The proposed algorithm performed slightly better than did the other methods using this dataset. The performance for each class is listed in Table 1. The proposed algorithm presented an AUROC value higher by 0.002 than the OCGAN value. In particular, the performance values for classes "2" and "8" were improved by 0.017 and 0.016, respectively, where the algorithms, including OCGAN, performed poorly because of the two classes

**TABLE 1.** Comparison of mean and variance of AUROC values for various methods (MNIST).

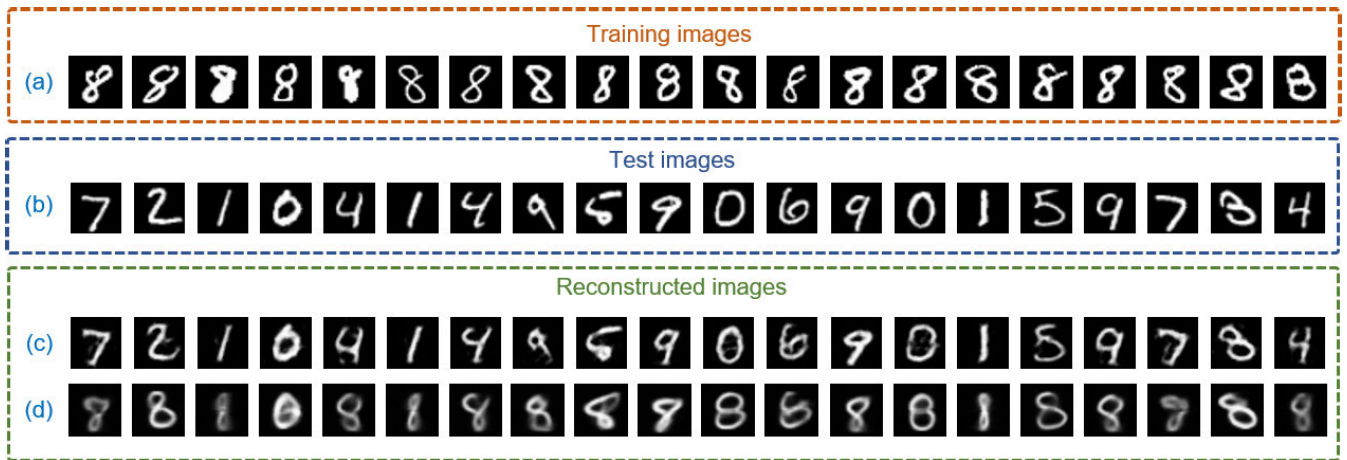| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | **MEAN** |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Autoencoder [1]** | 0.995 (0.001) | 0.999 (0.000) | 0.907 (0.006) | 0.945 (0.004) | 0.950 (0.006) | 0.961 (0.007) | 0.986 (0.002) | 0.966 (0.004) | 0.849 (0.012) | 0.966 (0.003) | 0.952 |
| **VAE [8]** | 0.997 (0.002) | 0.999 (0.000) | 0.935 (0.004) | 0.958 (0.006) | 0.973 (0.003) | 0.962 (0.005) | 0.991 (0.002) | 0.975 (0.004) | 0.921 (0.009) | 0.976 (0.006) | 0.966 |
| **AnoGAN [18]** | 0.963 (0.002) | 0.993 (0.001) | 0.849 (0.009) | 0.881 (0.007) | 0.894 (0.002) | 0.883 (0.005) | 0.943 (0.002) | 0.935 (0.003) | 0.839 (0.012) | 0.924 (0.003) | 0.910 |
| **GANomaly [20]** | 0.990 (0.004) | 0.999 (0.000) | 0.914 (0.012) | 0.936 (0.006) | 0.970 (0.002) | 0.966 (0.004) | 0.992 (0.001) | 0.977 (0.003) | 0.929 (0.010) | 0.979 (0.002) | 0.965 |
| **OCGAN [14]** | **0.998** **(0.001)** | 0.999 (0.000) | 0.947 (0.012) | **0.962** **(0.010)** | **0.972** **(0.004)** | **0.980** **(0.008)** | 0.991 (0.004) | **0.980** **(0.005)** | 0.938 (0.012) | 0.981 (0.003) | 0.975 |
| **Proposed** | **0.998** **(0.001)** | 0.999 (0.000) | **0.964** **(0.004)** | 0.958 (0.004) | 0.959 (0.003) | 0.978 (0.010) | **0.996** **(0.001)** | 0.976 (0.005) | **0.954** **(0.007)** | **0.984** **(0.002)** | **0.977** |



**FIGURE 4.** Illustration of the effect of decentralization loss on MNIST class "8": (a) images used in the training phase, (b) test images, (c) images reconstructed by the autoencoder [1], (d) images reconstructed by the proposed method with λ = 0.01, and (e) image reconstructed by the proposed method with λ = 1.

having similar characteristics as the number in other classes. It was noted that this improved performance supports the effectiveness of the proposed method, despite its simple architecture. As shown in Fig. 4, the autoencoder trained by numbers, such as "8" with complex shapes, generated well for numbers in classes other than "8." However, the model trained by the proposed method generated an image that resembled "8" for all input images, increasing the anomaly score for images that were not in class "8".

*b: FASHION MNIST*

The Fashion MNIST dataset comprises a set of grayscale images having a resolution of 28 × 28 and includes 10 classes of clothing and accessories. This dataset has larger between-class variation than does the MNIST dataset. Fig. 5 compares the reconstructed images of the two algorithms. Each algorithm was trained for the bag class of Fashion MNIST. The images in (b) were input images for testing, and the images in (c) and (d) were the output images from OCGAN and the proposed method, respectively, for the test images in (b). As shown in Fig. 5, the proposed method generated more variations of bag shapes than did OCGAN.

These results increased the difference between the input and output images, resulting in an anomaly-score increase for abnormal samples. Despite its simple structure, the proposed algorithm improved the AUROC by 0.006 compared with the state-of-the-art algorithm, OCGAN.

*c: MVTec*

The MVTec dataset consists of five texture classes and 10 object classes, with a total of 3,629 training images and 1,725 test images. The resolution of the images is 1024 × 1024. We resized these images to 256 × 256 for the experiments. The training data consist of one normal class, and the test data consist of one normal class and several abnormal classes. As shown in Fig. 6, unlike MNIST or Fashion MNIST, normal and abnormal images were not clearly clustered in MVTec. Therefore, existing methods of aggregating latent features showed low performance. The proposed method showed relatively good performance in this case, because it spread latent features to generate an image having the characteristics of a normal class at all times rather than simply expressing the distribution for a specific class. As shown in Fig. 3, the proposed method generated images
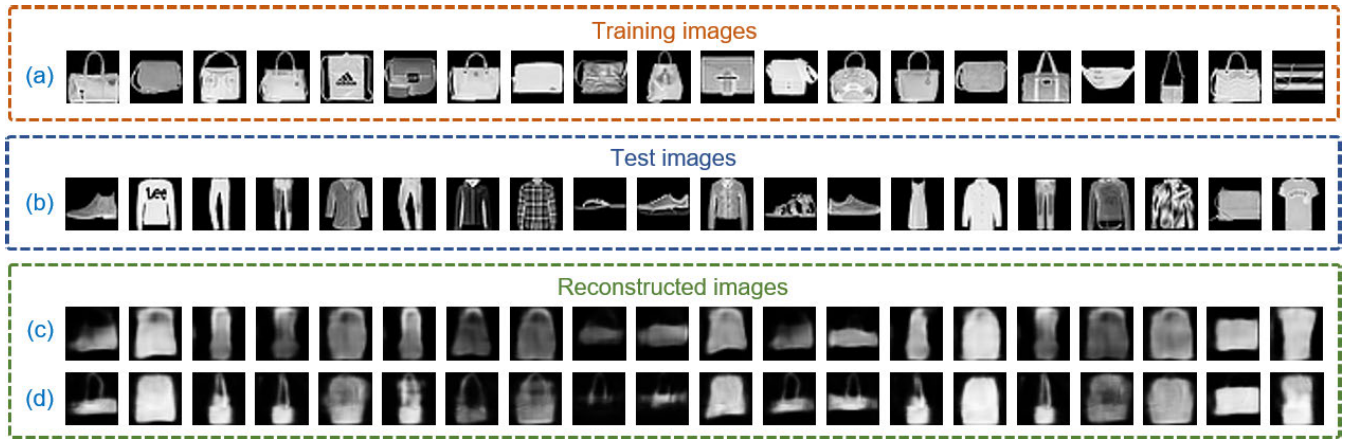
**FIGURE 5.** Illustration of the effect of decentralization loss on the Fashion MNIST class "bag": (a) sample images used in the training phase, (b) test images, (c) images reconstructed by OCGAN [14], and (d) images reconstructed by the proposed method.

**TABLE 2.** Comparison of mean and variance of AUROC values for various methods (fashion MNIST).

|  | T-shirts | Trouser | Pullover | Dress | Coat | Sandal | Shirt | Sneaker | Bag | Ankle boot | MEAN |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Autoencoder [1]** | 0.910 (0.002) | 0.980 (0.001) | 0.862 (0.003) | 0.919 (0.002) | 0.874 (0.010) | 0.908 (0.004) | 0.803 (0.008) | 0.978 (0.003) | 0.856 (0.001) | 0.965 (0.004) | 0.906 |
| **VAE [8]** | 0.896 (0.002) | 0.984 (0.003) | 0.875 (0.002) | 0.927 (0.004) | 0.899 (0.003) | 0.906 (0.005) | 0.830 (0.006) | 0.982 (0.001) | 0.859 (0.009) | 0.976 (0.001) | 0.913 |
| **AnoGAN [18]** | 0.824 (0.005 | 0.957 (0.007) | 0.812 (0.009) | 0.904 (0.006) | 0.861 (0.008) | 0.866 (0.011) | 0.798 (0.005) | 0.933 (0.002) | 0.801 (0.008) | 0.951 (0.004) | 0.871 |
| **GANomaly [20]** | 0.866 (0.013) | 0.978 (0.004) | 0.840 (0.008) | 0.930 (0.008) | 0.885 (0.005) | 0.918 (0.015) | 0.815 (0.006) | 0.977 (0.004) | 0.867 (0.014) | **0.980** **(0.005)** | 0.906 |
| **OCGAN [14]** | 0.901 (0.004) | 0.986 (0.002) | 0.887 (0.011) | **0.942** **(0.001)** | 0.901 (0.005) | **0.920** **(0.004)** | 0.835 (0.004) | 0.984 (0.001) | 0.876 (0.004) | **0.980** **(0.002)** | 0.921 |
| **Proposed** | **0.915** **(0.003)** | **0.990** **(0.001)** | **0.898** **(0.002)** | 0.938 (0.002) | **0.921** **(0.004)** | 0.912 (0.003) | **0.840** **(0.004)** | **0.989** **(0.001)** | **0.887** **(0.009)** | 0.977 (0.002) | **0.927** |

**TABLE 3.** Comparison of mean and variance of AUROC values for MVTEC (object) dataset.

|  | bottle | cable | capsule | hazelnut | metalnut | pill | screw | toothbrush | transistor | zipper | MEAN |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **OCGAN [14]** | 0.897 (0.002) | 0.798 (0.003) | 0.734 (0.002) | 0.721 (0.004) | **0.593** **(0.010)** | 0.797 (0.005) | 0.659 (0.008) | 0.847 (0.003) | 0.804 (0.009) | 0.712 (0.004) | 0.756 |
| **Proposed** | **0.944** **(0.005** | **0.835** **(0.007)** | **0.745** **(0.009)** | **0.738** **(0.006)** | 0.582 (0.008) | **0.818** **(0.011)** | **0.673** **(0.005)** | **0.856** **(0.002)** | **0.835** **(0.008)** | **0.746** **(0.004)** | 0.777 |

more similar to the normal images than did the images from any other method. This led to a large anomaly score for abnormal images. Furthermore, the proposed method exhibited better performance than did the OCGAN, and it improves performance in most classes.

In summary, we demonstrated the advantages of the proposed method through an experiment conducted on a total of three datasets. As shown in Tables 1–3, the decentralization loss significantly improved performance for all three datasets. This means that the proposed decentralization loss was effective for one-class anomaly detection, regardless of datasets. This is because, unlike the existing methods, the proposed method restores the characteristic image of the normal class, as shown in Figs. 3–5. Fig. 7 shows the distribution of the anomaly scores of two different AUROC values. It shows the effect of increasing the AUROC value on the distribution of the anomaly score. As the AUROC
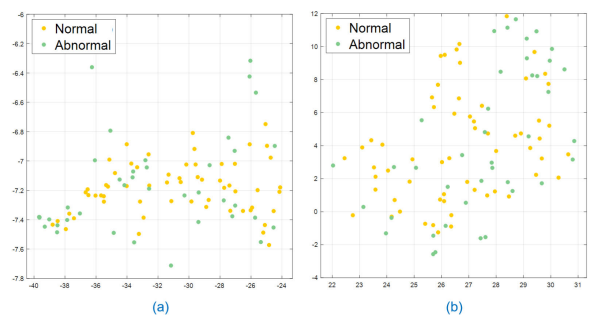


**FIGURE 6.** t-SNE [33] plot results for latent feature of MVTec Transistor test data. (a) latent features encoded by the autoencoder [1], (b) latent features encoded by the proposed method with λ = 0.01.

value increased, the difference between the anomaly scores of normal and abnormal samples increased. MNIST and Fashion MNIST are not datasets designed for anomaly detection;
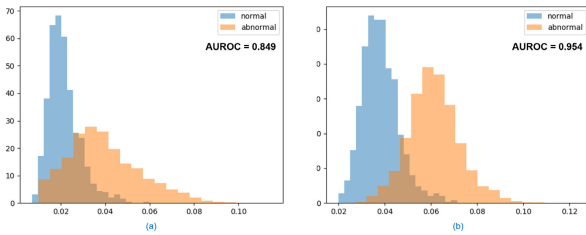
**FIGURE 7.** Illustration of histogram of anomaly scores for MNIST test data. (a) autoencoder [1](low AUROC value), (b) proposed (high AUROC value). Each model is trained on MNIST class "8".



**FIGURE 8.** Illustration of training time of proposed method and OCGAN [14]. Training epoch is set to 50.

the difference between normal and abnormal samples was distinct. As a result, most algorithms performed well for these datasets. The MVTec dataset, however, was designed for anomaly detection, and there are not many differences in the patterns or shapes of normal and abnormal images. Thus, most methods performed poorly, as shown in Fig. 3. OCGAN and the proposed algorithm performed better than did the existing algorithms using the approach that allowed the model to generate only a normal class image. However, unlike OCGAN, which trains a model based on randomly generated vectors, the proposed method optimizes the loss function based on the statistical properties of the latent vectors of the training data. Therefore, it is possible to find the optimal point where the MSE loss and decentralization loss are balanced. As a result, the proposed algorithm achieved a better result than did OCGAN for all datasets. The proposed algorithm achieved 0.01 better than the higher AUROC value on average for three datasets than that of OCGAN.

### 2) COMPLEXITY COMPARISON WITH OCGAN
OCGAN was proposed as a model for positioning latent features in all areas of a manifold space. However, OCGAN has a complex network architecture and complex learning schemes. As shown in Fig. 2, decentralization loss achieved the same objective while not requiring any additional, complex architecture or learning scheme during the training phase. Consequently, the proposed method was more efficient and easier to implement. The proposed method requires training on only half as many parameters as did OCGAN. As a result, as shown in Fig. 8, the time for required training was reduced by 48%. Moreover, rather than relying on random sampling, the proposed loss function sought to maximize the between-class variance more directly, which is extremely beneficial for one-class anomaly detection. As shown in Tables 1–3, the effectiveness of the proposed method was demonstrated by experimental results for various datasets.

### 3) PARAMETER SELECTION FOR $\lambda$
We also conducted experiments to evaluate how $\lambda$ influences latent-feature distribution and the reconstructed image. To estimate the effect of $\lambda$ and to show that the proposed method can provide a solution for various datasets, we considered multiple $\lambda$ values with multiple datasets. We conducted experiments on datasets comprising normal and abnormal
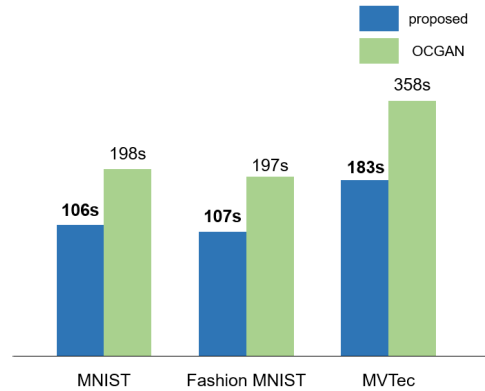
sample images containing different objects. We also identified the anomalies having similar structures in normal samples, such as in the MVTec dataset. Fig. 9 shows that different $\lambda$ values led to different deep-feature distributions: the larger the $\lambda$ value, the more the latent vector spread over the manifold space. With a properly selected $\lambda$, features were spread over a broader range on the manifold space, and reconstructed images were always similar to normal samples. However, a large $\lambda$ value, as shown in Fig. 9, causes the model to reconstruct the same image for all inputs. This degrades the performance of the algorithm, because, as the lambda value increases, the proportion of the MSE decreases in the overall loss function. Also, as shown in Table 4, the performance of the proposed algorithm varied with the $\lambda$ value up to 0.064. Therefore, the joint supervision and balancing of the benefits of the two components is crucial for one-class anomaly detection. In this paper, the optimal $\lambda$ value of 0.01 was determined through experiments.

**TABLE 4.** Experimental results for various $\lambda$ (AUROC).

|  | $\lambda = 1$ | 0.1 | 0.01 | 0.001 |
|---|---|---|---|---|
| **MNIST** | 0.932 | 0.948 | **0.977** | 0.967 |
| **Fashion MNIST** | 0.902 | 0.909 | **0.927** | 0.918 |
| **MVTec** | 0.718 | 0.737 | **0.777** | 0.766 |

### 4) ABLATION STUDY
An ablation study was conducted to show the effectiveness of the proposed regularization term and simplified regularization term. The experiment was conducted on a total of three datasets. As shown in Table 5, for all three datasets, the decentralization loss improved the performance. Additionally, the high performance was maintained, even for the model trained by (13). Although the model that applied (9) had a

**TABLE 5.** Experimental results for ablation study (AUROC).

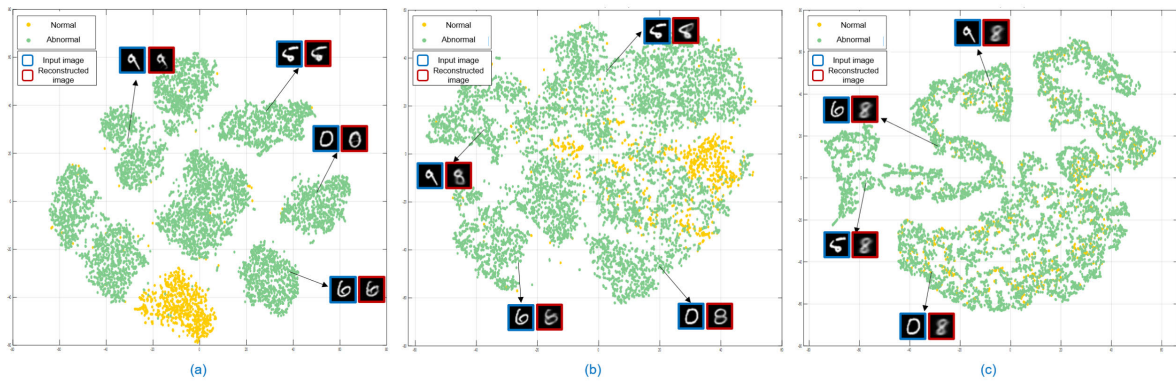|  | MNIST | Fashion MNIST | MVTec |
|---|---|---|---|
| **autoencoder** | 0.952 | 0.906 | 0.698 |
| **autoencoder with (9)** | 0.979 | 0.928 | 0.779 |
| **autoencoder with (13)** | 0.977 | 0.927 | 0.777 |

**FIGURE 9.** t-SNE [33] plot results for latent feature of MNIST test data. (a) latent features encoded by the autoencoder [1], (b) latent features encoded by the proposed method with λ = 0.01, and (c) latent features encoded by the proposed method with λ = 1. Each model is trained on MNIST class "8" as a normal class; The blue border images represents the input images, and the red border images means the output image for the input image.

slightly better performance, to emphasize the simple implementation, a strength of the proposed algorithm, the experimental results from (13) were used in this paper. When the proposed method was applied, the proposed method spread the latent-feature vectors over the broad areas, as shown in Fig. 8. As a result, the characteristic image of the normal class was restored for any input image, as shown in Figs. 3–5.

### 5) MULTI-MODAL DISTRIBUTION

Decentralization loss is a method designed upon unimodal data. In other words, by dispersing latent features based on one central tendency/mean point, an image have normal characteristics is displayed for any image. Additional experiments were conducted to evaluate the effect of decentralization loss on multi-modal data. For MNIST, two arbitrary classes were defined as normal, and the remaining seven classes were defined as abnormal. The equation used to optimize in the proposed method is (9). In the unimodal case, the latent feature spreads based on the mean that corresponds to the class by minimizing (9). However, in the multi-modal case, the latent feature spreads based on the center point of the mean of the two classes, and, when the latent features of the two classes are mixed, the MSE loss that corresponds to the former in (9) increases. Eventually, as shown in Fig. 10, latent features do not spread over a certain range around the mean of each class. Therefore, there is no negative effect, even in the multi-modal case. Rather, it can be extended to a multi-class classification network to provide the attention mode for specific class to distinguish from other classes with similar images via the proposed loss function, or it can be applied to various studies, such as open-set recognition and out-of-distribution [12], [13].

### E. LIMITATION AND FUTURE WORK

Minimizing the MSE and decentralization loss simultaneously induced the model to generate an image containing the common features of all normal data, increasing the reconstruction error for abnormal samples. Thus, the proposed
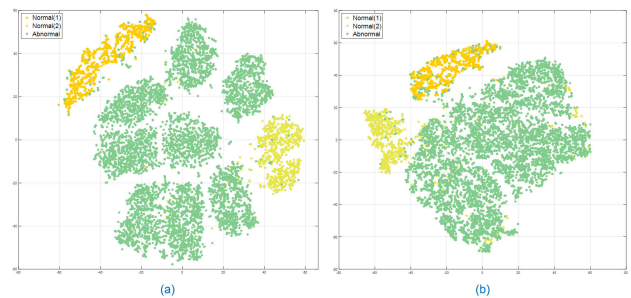


**FIGURE 10.** t-SNE [33] plot results for latent feature of MNIST test data. (a) latent features encoded by the autoencoder [1], (b) latent features encoded by the proposed method with λ = 0.01. Each model is trained on MNIST classes "2" and "3" as a normal class; The decentralization loss fail to spread latent vectors for learned classes to the entire space.

method is advantageous for distinguishing abnormal samples having a different tendency from the characteristics of normal data. The proposed loss function function tended to disappear at each detail of the image, owing to the fact that it was based on statistical properties, such as mean and variance. As shown in Fig. 11, the proposed algorithm exhibited poor performance on a dataset containing different details for each image. For example, this occurred for one in which the positions of the letters differed for all images. The green areas in Fig. 11 (d) indicated, for each column, the differences between the test image and the image reconstructed using the proposed method. In Fig. 11(a), the orange boxes show that the position of the letters in these normal images changed. Consequently, the letters were erased in the reconstructed image. Therefore, the anomaly score for abnormality and normality did not differ significantly, and the performance was degraded. This property is entirely different from those observed for images reconstructed using the existing methods in which all details of the images were reconstructed. Further work is essential to find the optimal point between conserving the individual images' details and encapsulating the characteristics of the class by introducing additional functions, such as perceptual loss [34] and U-Net [35].
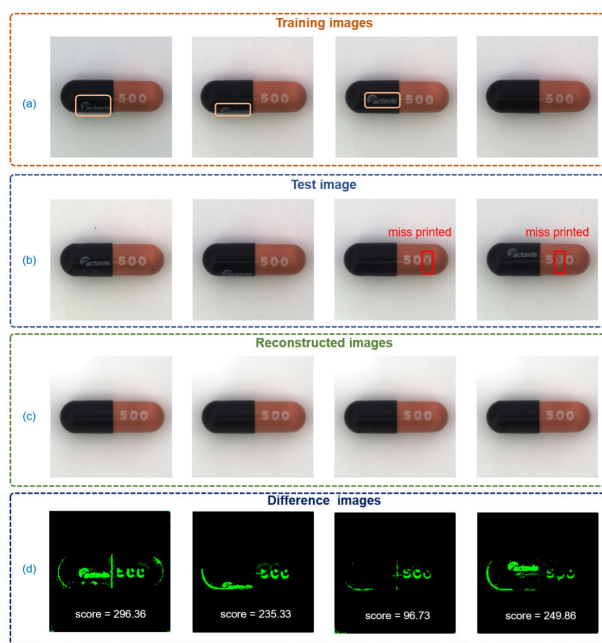
**FIGURE 11.** Illustration of the limitation of the proposed method. (a) Sample images used in the training phase; the position of the letters is indicated by an orange box in each one. (b) Test images; the two images on the left are normal, and the two images on the right are abnormal. (c) Images reconstructed by the proposed method; the input image is the image with the same number in (b). (d) The differences between the images in (b) and (c).

## V. CONCLUSION

In this paper, we proposed a novel objective function for one-class anomaly detection using a new approach for spreading vectors in manifold space. We designed the loss function based on the statistical properties of the latent vector, such as mean and variance, and restricted the candidate values of the latent vector, because the range of latent-vector values in each dataset is very different. Thus, the range of the proposed regularization term depends only on the size of the latent space. We experimentally found an optimal regularization parameter that could be applied regardless of the dataset's statistical characteristics. Additionally, compared with the state-of-the-art method, the proposed function achieved better performance by 0.002, 0.006, and 0.021 on the MNIST, Fashion MNIST, and MVTec datasets, respectively, despite its simple architecture. Furthermore, we achieved a performance improvement of approximately 1.2% for classes with which existing algorithms had difficulty. This is a significant improvement, given that the proposed method reduced training time by 48%. Moreover, although the proposed algorithm was designed for a single modal distribution, it had no side effects, even for a multi-modal distribution. Based on this, the model can be extended to various studies. As a result, the effectiveness of the proposed method was demonstrated for three datasets. However, some details were omitted in the reconstructed image for a detailed image dataset, because the proposed method utilized statistical properties. Further research will be conducted to find the balance between

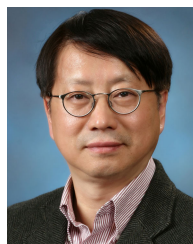conserving the details of each image and encapsulating the characteristics of all the data.

## REFERENCES

[1] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, and L. Bottou, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J. Mach. Learn. Res.*, vol. 11, no. 12, 2010.

[2] C. M. Bishop, *Pattern Recognition Machine Learning*. Cham, Switzerland: Springer, 2006.

[3] E. Eskin, A. Arnold, M. Prerau, L. Portnoy, and S. Stolfo, "A geometric framework for unsupervised anomaly detection: Detecting intrusions in unlabeled data," in *Application of Data Mining in Computer Security*. Norwell, MA, USA: Kluwer, 2002.

[4] M. Nicolau, "One-class classification for anomaly detection with kernel density estimation and genetic programming," in *Proc. Eur. Conf. Genetic Program.*, 2016, pp. 3–18.

[5] J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," Microsoft Res. (MSR), Bengaluru, India, Tech. Rep. MSR-T R-99-87, 1999.

[6] R. Ranjan and G. Sahoo, "A new clustering approach for anomaly intrusion detection," 2014, *arXiv:1404.2772*. [Online]. Available: http://arxiv.org/abs/1404.2772

[7] N. Görnitz, M. Braun, and M. Kloft, "Hidden Markov anomaly detection," in *Proc. Int. Conf. Mach. Learn.*, vol. 2015, pp. 1833–1842.

[8] D. P Kingma and M. Welling, "Auto-encoding variational bayes," 2013, *arXiv:1312.6114*. [Online]. Available: http://arxiv.org/abs/1312.6114

[9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. Adv. Neural Inf. Process. Syst.*, 2014, pp. 2672–2680.

[10] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, "Adversarial autoencoders," 2015, *arXiv:1511.05644*. [Online]. Available: http://arxiv.org/abs/1511.05644

[11] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mvtec ad–a comprehensive real-world dataset for unsupervised anomaly detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9592–9600.

[12] C.-K. Yeh, W.-C. Wu, W.-J. Ko, and Y.-C. F. Wang, "Learning deep latent space for multi-label classification," in *Proc. AAAI Conf. Artif. Intell.*, 2017, pp. 1–4.

[13] P. Oza and V. M. Patel, "C2AE: Class conditioned auto-encoder for open-set recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 2307–2316.

[14] P. Perera, R. Nallapati, and B. Xiang, "OCGAN: One-class novelty detection using GANs with constrained latent representations," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 2898–2906.

[15] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song, "Using self-supervised learning can improve model robustness and uncertainty," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 15637–15648.

[16] P. Bergmann, S. Löwe, M. Fauser, D. Sattlegger, and C. Steger, "Improving unsupervised defect segmentation by applying structural similarity to autoencoders," 2018, *arXiv:1807.02011*. [Online]. Available: http://arxiv.org/abs/1807.02011

[17] J. An and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," Seoul Nat. Univ., Seoul, South Korea, Tech. Rep. 2015-03, 2015.

[18] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs, "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *Proc. Int. Conf. Inf. Process. Med. Imag.*, 2017, pp. 146–157.

[19] H. Zenati, C. Sheng Foo, B. Lecouat, G. Manek, and V. R. Chandrasekhar, "Efficient GAN-based anomaly detection," 2018, *arXiv:1802.06222*. [Online]. Available: http://arxiv.org/abs/1802.06222

[20] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon, "Ganomaly: Semi-supervised anomaly detection via adversarial training," in *Proc. Asian Conf. Comput. Vis.*, 2018, pp. 622–637.

[21] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli, "Adversarially learned one-class classifier for novelty detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 3379–3388.

[22] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon, "Skip-GANomaly: Skip connected and adversarially trained encoder-decoder anomaly detection," in *Proc. Int. Joint Conf. Neural Netw. (IJCNN)*, Jul. 2019, pp. 1–8.

[23] C. Phuc Ngo, A. Aristo Winarto, C. Kou Khor Li, S. Park, F. Akram, and H. Kuan Lee, "Fence GAN: Towards better anomaly detection," 2019, *arXiv:1904.01209*. [Online]. Available: http://arxiv.org/abs/1904.01209

[24] T. Popoviciu, "Sur leséquations algébriques ayant toutes leurs racines réelles," *Mathematica*, vol. 9, pp. 129–145, Oct. 1935.

[25] R. Bhatia and C. Davis, "A better bound on the variance," *Amer. Math. Monthly*, vol. 107, no. 4, pp. 353–357, Apr. 2000.

[26] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 499–515.

[27] T. Pang, C. Du, and J. Zhu, "Max-mahalanobis linear discriminant analysis networks," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 4016–4025.

[28] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv:1409.1556*. [Online]. Available: http://arxiv.org/abs/1409.1556

[29] T. Pang, K. Xu, Y. Dong, C. Du, N. Chen, and J. Zhu, "Rethinking softmax cross-entropy loss for adversarial robustness," in *Proc. Int. Conf. Learn. Represent.*, 2019, pp. 1–19.

[30] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," 2013, *arXiv:1312.6199*. [Online]. Available: http://arxiv.org/abs/1312.6199

[31] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 770–778.

[32] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 630–645.

[33] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, pp. 2579–2605, Nov. 2008.

[34] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 694–711.

[35] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assist. Intervent*, 2015, pp. 234–241.

**EUNGI HONG** (Graduate Student Member, IEEE) received the B.S. degree in electrical engineering from Yonsei University, Seoul, South Korea, in 2013, where he is currently pursuing the Ph.D. degree in electrical and electronic engineering. His research interests include image classification, open-set recognition, and anomaly detection.

**YOONSIK CHOE** (Senior Member, IEEE) received the B.S. degree in electrical engineering from Yonsei University, Seoul, South Korea, in 1979, the M.S.E.E. degree in systems engineering from Case Western Reserve University, Cleveland, OH, USA, in 1984, the M.S. degree in electrical engineering from Pennsylvania State University, State College, PA, USA, in 1987, and the Ph.D. degree in electrical engineering from Purdue University, West Lafayette, IN, USA, in 1990. From 1990 to 1993, he was a Principal Research Staff with the Industrial Electronics Research Center, Hyundai Electronics Company Ltd. Since 1993, he has been with the Department of Electrical and Electronic Engineering, Yonsei University, where he is currently a Professor. His research interests include video coding, video communication, statistical signal processing, and digital image processing.

● ● ●