

Deep Unsupervised Image Anomaly Detection: An Information Theoretic Framework

Fei Ye, Huangjie Zheng, Chaoqin Huang, Ya Zhang

Abstract—Surrogate task based methods have recently shown great promise for unsupervised image anomaly detection. However, there is no guarantee that the surrogate tasks share the consistent optimization direction with anomaly detection. In this paper, we return to a direct objective function for anomaly detection with information theory, which maximizes the distance between normal and anomalous data in terms of the joint distribution of images and their representation. Unfortunately, this objective function is not directly optimizable under the unsupervised setting where no anomalous data is provided during training. Through mathematical analysis of the above objective function, we manage to decompose it into four components. In order to optimize in an unsupervised fashion, we show that, under the assumption that distribution of the normal and anomalous data are separable in the latent space, its lower bound can be considered as a function which weights the trade-off between mutual information and entropy. This objective function is able to explain why the surrogate task based methods are effective for anomaly detection and further point out the potential direction of improvement. Based on this object function we introduce a novel information theoretic framework for unsupervised image anomaly detection. Extensive experiments have demonstrated that the proposed framework significantly outperforms several state-of-the-arts on multiple benchmark data sets. Source code will be made publicly available.

Index Terms—Anomaly detection, information theoretic framework, mutual information, entropy

I. INTRODUCTION

ANOMALY detection, aiming to find patterns that do not conform to expected behavior, has been broadly applicable in medical diagnosis, credit card fraud detection, sensor network fault detection and numerous other fields [1]. With the recent advance of the convolution neural network, anomaly detection for image data has received significant attention, which is challenging due to high-dimension and complex texture of images. As anomalous data is diverse and inexhaustible, anomaly detection is usually considered under the setting where the training data contains only the “normal” class, *i.e.*, the anomalous data is not available during training.

In this paper, we explore unsupervised anomaly detection for images. Anomaly detection aims to learn effective feature representation that leads to the best separation of anomalous data from normal data. However, under the unsupervised setting, only the normal data is available during the model fitting.

It seems straightforward to make assumptions on the distribution of anomalous data and model unsupervised anomaly detection as a one-class classification problem. Several methods have explored in this direction. Scholkopf et al. [2] and Oza et al. [3] assumed the anomalous distribution to be the origin and the zero centered Gaussian distribution accordingly. Ruff et al. [4] assumed that the anomalous distribution to be a uniform distribution in latent space. However, the simple assumption of the anomalous distribution can not force the network to extract effective features. Making it inevitable to rely on the pre-trained model, resulting in a two-step optimizing process.

Other than directly learning to separate anomalous and normal data, surrogate based approaches attempt to first learn feature representation with an unsupervised alternative objective function, and then assume that model will lead to poor performance to the anomalous data, which are not exposed to the training, so that they can be distinguished from the normal data. Two types of surrogate tasks are typically adopted: reconstruction [5]–[8] and self-labeling [9], [10]. For self-labeling based method, Golan et al. [9] applied dozens of image geometric transforms and created a self-labeled dataset for transformation classification. Wang et al. [10] introduced more self-label method like patch re-arranging and irregular affine transformations. Recently, a new surrogate task: restoration [11] is introduced which assumes that through restoring the erased information, the model is effectively forced to learn what is erased, and the feature embedding can thus be controlled by the corresponding information erasing. While surrogate approaches have become the mainstream method for anomaly detection in recent years and have shown promising results, it is hard to ensure that the surrogate tasks share the consistent optimization direction with anomaly detection.

In this paper, we return to a direct objective function for anomaly detection, which maximizes the distance between normal and anomalous data in terms of the joint distribution for image and feature representation. We decompose the above objective function into the following four components through mathematical analysis: *Mutual information* between image space and latent space of normal data, *Entropy* of normal data in latent space; *Expectation of cross-entropy* between normal and anomalous samples in latent space; *Distribution distance* between normal and anomalous samples in image space. The first two components can be calculated with normal data only. The fourth component is a constant once normal data is selected. Only optimizing the third term requires anomalous data. To optimize the objective function under the unsupervised setting, we investigate the condition to bypass the third term and get a lower bound on the objective function which can

F. Ye, C. Huang and Y. Zhang are with the Cooperative Medianet Innovation Center and the Shanghai Key Laboratory of Multimedia Processing and Transmissions, Shanghai Jiao Tong University, Shanghai 200240, China. (E-mail: {yf3310, huangchaoqin, ya_zhang}@sjtu.edu.cn). H. Zheng is with Department of statistics and data science, University of Texas at Austin, Austin, 78751 TX, USA. (E-mail: huangjie.zheng@utexas.edu).

be considered as a trade-off between the mutual information and the entropy. To our best knowledge, this is the first end-to-end framework to optimize anomaly detection directly. We provide a specific method based on the framework and further show that the lower-bound objective function can be linked to several previous studies such as the reconstruction-based method and the classic one-class classification method SVDD [4]. As most approaches focus on only one term (mutual information or entropy), the proposed objective function can not only fill the lack of theory for many of the existing anomaly detection method but also point out the potential direction of improvement. To our best knowledge, this is the first anomaly detection objective function that can be end-to-end optimized under the unsupervised setting.

To validate the effectiveness of our lower bound objective function, we conduct extensive experiments on several benchmark datasets, including MNIST [12], Fashion-MNIST [13], CIFAR-10 [14], CIFAR-100 [14] and ImageNet [15]. Our experimental results have shown that the proposed method outperforms several state-of-the-art methods in terms of model accuracy and model stability to a large extent. The main contributions of the paper are summarized as follows:

II. RELATED WORKS

A. Anomaly Detection

For anomaly detection on images and videos, a large variety of methods have been developed in recent years [1], [16]–[23]. As anomalous data is diverse and inexhaustible, anomaly detection is usually considered under the setting where the training dataset contains only “normal” data, i.e. the anomalous data is not available during training. The main challenges are two-fold, the first one is how to extract effective features, the second one is how to separate anomalous data in latent space.

All the previous approaches can be broadly classified into three categories, statistic based, surrogate based and one-class classification based approaches.

Statistic based approaches [24]–[26] assume that anomalous data will be mapped to different statistical representations that are far from the distribution of normal data. The features are typically extracted by some shallow encoders like Principal Component Analysis (PCA), kernel PCA, or Robust PCA [27].

Surrogate based approaches assume that the anomalous data will yield in different embedding from normal data and lead to poor performance which can be utilized as criteria to define anomaly. It manages to tackle the first challenge through unsupervised learning with an alternative objective function other than optimizing anomaly detection directly. Three main frameworks, different in the employed supervision, are typically adopted: reconstruction-based, self-labeling-based, and outlier exposure-based frameworks.

Reconstruction-based frameworks take input image as supervision and assume that anomalous data have larger reconstruction error. The advantage is that supervision can be easily obtained. The main challenge is to find a more effective loss function to replace the typically adopted pixel-wised MSE loss, which is indicated ineffective to force the model to extract discriminate features [28], [29]. Adversarial training

is leveraged to optimize the pixel-wised MSE loss, through adding a discriminator after autoencoders to judge whether its original or reconstructed image [5], [6]. Akcay et al. [7] adds an extra encoder after autoencoders and enclosing the distance between the embedding. Perera et al. [8] applied two adversarial discriminators and a classifier on a denoising autoencoder. By adding constraint and forcing each randomly drawn latent code to reconstruct examples like the normal data, it obtained high reconstruction errors for the anomalous data.

Self-label based frameworks, which take the artificial label as supervision and assume that the labels can not be predicted properly for anomalous data, has recently received significant attention. This framework, as decoder-free, can be a benefit in the advanced classification network which is proved to be more effective in extracting discriminate features. The main challenge is how to define meaningful labels. Golan et al. [9] applied dozens of image geometric transforms and created a self-labeled dataset for transformation classification. Wang et al. [10] introduced a more self-label method like patch re-arranging and irregular affine transformations.

Outlier exposure based frameworks take an auxiliary dataset entirely disjoint from test-time data as supervision and thus teach the network better representations for anomaly detection. Hendrycks et al. [30] introduced extra data to build a multi-class classification task. The experiment revealed that even though the extra data was in limited quantities and weakly correlated to the normal data, the learned hyperplane was still effective in separating normal data.

Restoration based framework, a new framework introduced by Ye et al. [11], which assumed that through restoring the erased information the model will be effectively forced to learn what is erased and how to restore it. Thus feature embedding can be controlled by the corresponding information erasing.

One-class classification based approaches tackle the second challenge by making assumptions on the distribution of anomalous data, thus change the anomaly detection into a supervised binary classification problem. Explicitly, Scholkopf et al. [2] and Oza et al. [3] assumed the anomalous distribution to be the origin and the zero centered Gaussian distribution in latent space accordingly. Implicitly, Ruff et al. [4] assumed that the anomalous distribution to be a uniform distribution in latent space. Despite these approaches [2]–[4], [31] could directly optimize the anomaly detection based objective function, the disadvantage is also obvious: these approaches have to rely on a pre-defined or pre-trained feature extractor, as its objective function and simple assumption on anomalous distribution can not force the network to extract effective feature.

B. Unsupervised Learning by Maximizing Mutual Information

Mutual information, generally employed to measure the correlation between two random variables, has recently received extensive attention in unsupervised learning. Mutual information is notoriously difficult to compute, particularly in continuous and high-dimensional settings [32]. To tackle this problem, Belghazi et al. [33] employ deep neural networks to effectively compute the mutual information between high dimensional input/output pairs. Hjelm et al. [32] utilize an 1x1

convolutional discriminator to compute the mutual information between a global summary feature vector and a collection of local feature vectors. Oord et al. [34] introduce a new NCE based loss called InfoNCE to maximize mutual information. Chen et.al [35] propose to maximize the agreement between two augmentation results from the same raw data and investigate properties of contrastive learning. Bachman et.al [36] maximize mutual information between features extracted from multiple views of a shared context. He et al. [37] utilize a dynamic dictionary to decouple the dictionary size from the mini-batch size, allowing the model to benefit from a large sampling size.

III. METHODS

A. Problem Formulation

Let \mathcal{X} and \mathcal{Z} be the domain of image data and representation data. Denote \mathcal{X}_n as the normal data and \mathcal{X}_a be the anomalous data, where $\mathcal{X}_n = \{\mathbf{x}_n : \mathbf{x}_n \sim p_n(\mathbf{x})\}$, $\mathcal{X}_a = \{\mathbf{x}_a : \mathbf{x}_a \sim p_a(\mathbf{x})\}$. Let the corresponding representation of \mathcal{X}_n and \mathcal{X}_a to be \mathcal{Z}_n and \mathcal{Z}_a accordingly, where $p_n(\mathbf{z})$ and $p_a(\mathbf{z})$ are the marginal distribution of the latent representations, $p_n(\mathbf{x}, \mathbf{z})$ and $p_a(\mathbf{x}, \mathbf{z})$ are the joint distribution between image space and the latent space.

To distinguish the normal data \mathbf{x}_n and anomalous data \mathbf{x}_a , due to the curse of dimensionality, anomaly detection in high-dimensional image space is not favorable in general cases. A common choice is to define a continuous and (almost everywhere) differentiable parametric function, $E_\theta : \mathcal{X} \rightarrow \mathcal{Z}$ with parameters θ (e.g., a neural network) that encodes the data to the low-dimensional latent space \mathcal{Z} , where the anomalous data can be detected more easily. Thus a straight-forward objective function may be maximizes the conditional distributions between normal and anomalous data in latent space. We take the commonly used KL divergence as measurement metric, and objective function can be formulate as:

$$\max_{\theta} \text{KL} [p_n(\mathbf{z}|\mathbf{x}) \parallel p_a(\mathbf{z}|\mathbf{x})]. \quad (1)$$

In order to make the object function in Eq. 1 more complete, in which the distribution distance in image space is taken into consideration simultaneously, we introduce a novel anomaly detection based objective function, which maximizes the distance between normal and anomalous data in terms of the joint distribution for image and feature representation. The objective function can be formulate as:

$$\max_{\theta} \text{KL} [p_n(\mathbf{x}, \mathbf{z}) \parallel p_a(\mathbf{x}, \mathbf{z})]. \quad (2)$$

When maximizing Eq. 2, all marginal distributions and all conditional distributions also match this maximization:

Remark 1. If $\text{KL} [p_n(\mathbf{x}, \mathbf{z}) \parallel p_a(\mathbf{x}, \mathbf{z})]$ is maximized, then it is equivalent that $\text{KL} [p_n(\mathbf{x}) \parallel p_a(\mathbf{x})]$ and $\text{KL} [p_n(\mathbf{z}|\mathbf{x}) \parallel p_a(\mathbf{z}|\mathbf{x})]$ are maximized.

Proof. The KL divergence for the joint distributions can be decomposed with chain rule [38]:

$$\begin{aligned} \text{KL} [p_n(\mathbf{x}, \mathbf{z}) \parallel p_a(\mathbf{x}, \mathbf{z})] &= \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{x}, \mathbf{z})}{p_a(\mathbf{x}, \mathbf{z})} \right] \\ &= \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{x})}{p_a(\mathbf{x})} + \log \frac{p_n(\mathbf{z}|\mathbf{x})}{p_a(\mathbf{z}|\mathbf{x})} \right] \\ &= \text{KL} [p_n(\mathbf{x}) \parallel p_a(\mathbf{x})] + \mathbb{E}_{p_n(\mathbf{x})} [\text{KL} [p_n(\mathbf{z}|\mathbf{x}) \parallel p_a(\mathbf{z}|\mathbf{x})]]. \end{aligned} \quad (3)$$

To maximize the KL divergence for the joint distributions is equivalent to maximize the KL divergence for both marginal and conditional distributions [39]. \square

B. Objective Function Decomposition

Remark 1 shows the equivalence of our objective to the common knowledge of anomaly detection. In this section we present that this objective yields a lower bound which can be optimized without anomalous data. It is a great challenge to optimize the objective function in Eq. 2 directly, since $p_a(\mathbf{x}, \mathbf{z})$ is not accessible. To tackle this challenge, firstly we decompose the objective function into four components, as introduced in the Proposition 1. Then, as shown in Proposition 2, we investigate the condition where we can bypass the component correlated with anomalous data and thus reformulated the the objective function in Eq. 2 into a lower-bound in Eq. 10 which can be optimized only with normal data. Finally, as shown in Corollary 2.1, the lower-bound can be optimized with a regularized Lagrange multiplier and get Eq. 13 as the final objective function. The Proposition 1 is introduced as follows:

Proposition 1. Let $I_n(\mathbf{x}, \mathbf{z})$, $H_n(\mathbf{z})$, $H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))$, $\text{KL} [p_n(\mathbf{x}) \parallel p_a(\mathbf{x})]$ denote the mutual information between \mathbf{x} and \mathbf{z} for normal data, the entropy of \mathbf{z} for normal data, the cross entropy between $p_n(\mathbf{z}|\mathbf{x})$ and $p_a(\mathbf{z}|\mathbf{x})$ and KL divergency between $p_n(\mathbf{x})$ and $p_a(\mathbf{x})$, respectively. Note that when the dataset is given, *i.e.* $p_n(\mathbf{x})$ and $p_a(\mathbf{x})$ are fixed (even though anomalous data is unknown). The objective function can be reformulated as:

$$\begin{aligned} &\max_{\theta} \text{KL} [p_n(\mathbf{x}, \mathbf{z}) \parallel p_a(\mathbf{x}, \mathbf{z})] \\ &= \max_{\theta} \{I_n(\mathbf{x}, \mathbf{z}) - H_n(\mathbf{z}) + \mathbb{E}_{p_n(\mathbf{x})} [H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))] \\ &\quad + \text{KL} [p_n(\mathbf{x}) \parallel p_a(\mathbf{x})]\}. \end{aligned} \quad (4)$$

Proof. The objective function in Eq. 2 can be reformulated as:

$$\begin{aligned} &\max_{\theta} \text{KL} [p_n(\mathbf{x}, \mathbf{z}) \parallel p_a(\mathbf{x}, \mathbf{z})] \\ &= \max_{\theta} \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{x}, \mathbf{z})}{p_a(\mathbf{x}, \mathbf{z})} \right] \\ &= \max_{\theta} \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{z}|\mathbf{x}) \cdot p_n(\mathbf{x})}{p_a(\mathbf{z}|\mathbf{x}) \cdot p_a(\mathbf{x})} \right] \\ &= \max_{\theta} \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{z}|\mathbf{x}) \cdot p_n(\mathbf{x}) \cdot p_n(\mathbf{z})}{p_a(\mathbf{z}|\mathbf{x}) \cdot p_a(\mathbf{x}) \cdot p_n(\mathbf{z})} \right] \\ &= \max_{\theta} \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \left(\frac{p_n(\mathbf{z}|\mathbf{x})}{p_n(\mathbf{z})} \cdot p_n(\mathbf{z}) \cdot \frac{1}{p_a(\mathbf{z}|\mathbf{x})} \cdot \frac{p_n(\mathbf{x})}{p_a(\mathbf{x})} \right) \right]. \end{aligned}$$

The above formula can be then decomposed into 4 components. The 1st component refers to the mutual information between the data sample \mathbf{x} and its representation \mathbf{z} for normal data:

$$\begin{aligned} \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{z}|\mathbf{x})}{p_n(\mathbf{z})} \right] &= \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{z}|\mathbf{x}) \cdot p_n(\mathbf{x})}{p_n(\mathbf{x}) \cdot p_n(\mathbf{z})} \right] \\ &= \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{x}, \mathbf{z})}{p_n(\mathbf{x}) \cdot p_n(\mathbf{z})} \right] = I_n(\mathbf{x}, \mathbf{z}). \end{aligned} \quad (5)$$

The 2nd component is the negative entropy of \mathbf{z} w.r.t p_n :

$$\mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} [\log p_n(\mathbf{z})] = -\mathbb{E}_{p_n(\mathbf{z})} \left[\log \frac{1}{p_n(\mathbf{z})} \right] = -H_n(\mathbf{z}). \quad (6)$$

The 3rd component is the expected value of the cross entropy between the conditional distributions $p_a(\mathbf{z}|\mathbf{x})$ and $p_n(\mathbf{z}|\mathbf{x})$:

$$\begin{aligned} \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{1}{p_a(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{p_n(\mathbf{x})} \mathbb{E}_{p_n(\mathbf{z}|\mathbf{x})} [-\log p_a(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{p_n(\mathbf{x})} [H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))]. \end{aligned} \quad (7)$$

With $p_n(\mathbf{x})$ and $p_a(\mathbf{x})$ fixed, the 4th component is a constant:

$$\mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{x})}{p_a(\mathbf{x})} \right] = \text{KL}[p_n(\mathbf{x})||p_a(\mathbf{x})] = C. \quad (8)$$

Thus the objective function in Eq. 2 can be reformulated as:

$$\begin{aligned} \max_{\theta} \text{KL}[p_n(\mathbf{x}, \mathbf{z})||p_a(\mathbf{x}, \mathbf{z})] \\ = \max_{\theta} \{I_n(\mathbf{x}, \mathbf{z}) - H_n(\mathbf{z}) + \mathbb{E}_{p_n(\mathbf{x})} [H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))] \\ + KL[p_n(\mathbf{x})||p_a(\mathbf{x})]\} \end{aligned} \quad (9)$$

□

The decomposed objective function in Eq. 4 is an essential general formula for optimizing anomaly detection, which combines unsupervised learning (1st and 2nd components) and supervised learning (3rd component).

As the 1st and 2nd components can be trained through an unsupervised fashion and force the encoder to extract effective features, the demand for anomalous data is greatly reduced.

To deal with unsupervised setting where anomalous data is complete absence during training, we seek to get rid of the 3rd component, which partially relies on anomalous data.

Proposition 2. If $p_n(\mathbf{x}, \mathbf{z})$ and $p_a(\mathbf{x}, \mathbf{z})$ have a certain distance such that for most samples $\mathbf{x}, \mathbf{z} \sim p_n(\mathbf{x}, \mathbf{z})$ the evaluated density $p_a(\mathbf{z}|\mathbf{x})$ is small enough, such that $p_a(\mathbf{z}|\mathbf{x}) \leq p_n(\mathbf{z})$ and $p_a(\mathbf{z}|\mathbf{x}) \leq 1$ almost everywhere,

then we can derive a lower bound of Objective function for Eq. 4:

$$\max_{\theta} \{I_n(\mathbf{x}, \mathbf{z}) - H_n(\mathbf{z})\}. \quad (10)$$

Proof. With the assumption that the evaluated density $p_a(\mathbf{z}|\mathbf{x})$ is small enough for most of the samples $\mathbf{x}, \mathbf{z} \sim p_n(\mathbf{x}, \mathbf{z})$. This ensures the non-negativity of $\mathbb{E}_{p_n(\mathbf{x})} [H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))]$ with the following inequality:

$$\begin{aligned} &\inf \mathbb{E}_{p_n(\mathbf{x})} [H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))] \\ &= \inf \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} [-\log p_a(\mathbf{z}|\mathbf{x})] \\ &\geq \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} [\inf (-\log p_a(\mathbf{z}|\mathbf{x}))] \\ &\geq 0. \end{aligned} \quad (11)$$

Moreover, the 4th component in Eq. 4 is a constant greater than zero. Then we have a lower bound to Eq. 4:

$$\text{KL}[p_n(\mathbf{x}, \mathbf{z})||p_a(\mathbf{x}, \mathbf{z})] \geq I_n(\mathbf{x}, \mathbf{z}) - H_n(\mathbf{z}). \quad (12)$$

□

The objective function in Eq. 10 is vital as it reveals a general formula for optimizing anomaly detection in an unsupervised fashion. Note that the assumption in Proposition 2 is appropriate in the anomaly detection tasks. Since we often have access to data samples instead of the true data distribution, considering the empirical distribution $p_n(\mathbf{z}) = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{z}_i}$; $\mathcal{Z}_n = \{\mathbf{z}_i\}_{i=1}^N$, we always have $p_n(\mathbf{z}) \leq 1$ and we have the evaluated density $p_a(\mathbf{z}_i|\mathbf{x}_i) \leq 1$. Moreover, with this assumption, we can further ensure the objective can be optimized with a regularized Lagrange multiplier.

Corollary 2.1. With the assumption in proposition 2.1, the lower bound shown in Equation 10 yields an objective function to be optimized with a entropy-regularized Lagrange multiplier:

$$\max_{\theta} \{I(\mathbf{x}, \mathbf{z}) - \beta \cdot H(\mathbf{z})\}, \quad (13)$$

where $\beta \geq 0$ is a positive coefficient for the weight of the entropy regularization.

Proof. To show Eq. (13) is a lower bound of Eq. (4), we only need to show the lower bound holds when $\beta = 0$:

$$\text{Eq. (4)} - \text{Eq. (13)} = \mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{z})}{p_a(\mathbf{z}|\mathbf{x})} \right] + \text{KL}[p_n(\mathbf{x})||p_a(\mathbf{x})].$$

$\mathbb{E}_{p_n(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_n(\mathbf{z})}{p_a(\mathbf{z}|\mathbf{x})} \right] \geq 0$ as the assumption $\frac{p_n(\mathbf{z})}{p_a(\mathbf{z}|\mathbf{x})} \geq 1$ in Proposition 2 while the second term is a positive constant, thus Eq. (4) – Eq. (13) ≥ 0 and completes the proof. □

C. Optimization

In this section, we extend our general lower-bound objective to concrete loss functions. *Start from here, as all the terms in Eq. 10 and 13 refer to normal data only, we omit the subscript ‘n’ and ‘a’ to specify normal and anomalous.* Moreover, we present the probabilistic expressions in the empirical distribution since we work with samples from the data distribution.

A challenge to maximize the objective function is that both mutual information and entropy are not always tractable. To maximize the mutual information between \mathbf{x} and \mathbf{z} , we apply the Contrastive Predictive Coding lower bound with Noise Contrastive Estimation [40]:

$$I(\mathbf{x}; \mathbf{z}) \geq \mathbb{E} \left[\frac{1}{K} \sum_{i=1}^K \log \frac{e^{f(x_i, z_i)}}{\frac{1}{K} \sum_{j=1}^K e^{f(x_i, z_j)}} \right] \triangleq I_{\text{NCE}}, \quad (14)$$

where f is the critic function that maps the inputs into a value in \mathbb{R} , which can be modeled with various methods, such as a similarity function or a neural discriminator. Here, x_i and $z_i = E_{\theta}(x_i)$ are called positive pairs; while x_i and $z_j = E_{\theta}(x_j \mathbb{I}_{[j \neq i]})$, where $\mathbb{I}_{[j \neq i]} \in \{0, 1\}$ is an indicator function evaluating to 1 iff $j \neq i$, are called negative pairs. Empirically, two independently randomly augmented versions of the same

Algorithm 1: Training Pseudocode for base model

Input: batch size N , similarity function sim , structure of model E_θ , entropy weight β , constant c_1 represents the size of z_i , constant c_2 is set to 20 for all experiment

- 1 **for** sampled minibatch $\{x_k\}_{k=1}^N$ **do**
- 2 **for all** $k \in \{1, \dots, N\}$ **do**
- 3 randomly draw two augmentation functions
 $t, t' \sim \mathcal{T}$
- 4 # Augmentation
- 5 $\tilde{x}_{2k-1} = t(x_k)$
- 6 $\tilde{x}_{2k} = t'(x_k)$
- 7 # Extract features
- 8 $z_{2k-1} = E_\theta(\tilde{x}_{2k-1})$
- 9 $z_{2k} = E_\theta(\tilde{x}_{2k})$
- 10 **for all** $i \in \{1, \dots, 2N\}$ and $j \in \{1, \dots, 2N\}$ **do**
- 11 # Similarity
- 12 $s_{i,j} = sim(z_i, z_j) = z_i^\top \cdot z_j$
- 13 $s'_{i,j} = c_2 \cdot \tanh\left(\frac{s_{i,j}}{c_1 \cdot c_2}\right)$
- 14 **define** $S = \{s'_{i,j}, i, j \in 2N\}$, $s_{max} = \max\{S\}$
- 15 $S_{shift} = S - s_{max}$
- 16 **define** $S_{shift} = \{\hat{s}_{i,j}, i, j \in 2N\}$
- 17 **define** $\ell(i, j) = -\log \frac{\exp(\hat{s}_{i,j})}{\sum_{k=1}^{2N} \mathbb{I}_{[k \neq i]} \exp(\hat{s}_{i,k})}$
- 18 $\mathcal{L}_{nce} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$
- 19 $\mathcal{L}_{entropy} = \frac{1}{2N} \sum_{k=1}^{2N} \|z_k\|_q$
- 20 $\mathcal{L} = \mathcal{L}_{nce} + \beta \cdot \mathcal{L}_{entropy}$
- 21 update network E_θ
- 22 **return** encoder network E_θ

sample, e.g. an image and its rotated view, are often used as positive pairs [35], [36].

To minimize the entropy, we seek the lower bound for the negative entropy as:

$$\begin{aligned}
 -H(\mathbf{z}) &= \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log p(\mathbf{z})] \\
 &= \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log p(\mathbf{z}) - \log r(\mathbf{z})] + \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log r(\mathbf{z})] \\
 &= \text{KL}[p_n(\mathbf{z})||r(\mathbf{z})] + \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log r(\mathbf{z})] \\
 &\geq \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log r(\mathbf{z})],
 \end{aligned} \tag{15}$$

where $r(\mathbf{z})$ is a reference distribution and $\text{KL}[p_n(\mathbf{z})||r(\mathbf{z})] \geq 0$. One common choice is to let $r(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$. Then, this lower bound is proportional to the L2 norm:

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log r(\mathbf{z})] \propto -\frac{1}{N} \sum_{i=1}^N \|z_i\|_2^2, \tag{16}$$

where $\|\cdot\|_2$ denotes the Frobenius norm when $p = 2$ and $z_i = E_\theta(x_i)$ denotes the latent feature of the i^{th} normal sample. Eq. 16 also provides us with a geometrical interpretation in the latent space. In this view, we can also generalize the Frobenius case into other orders, such as $p = 1$: Another choice is thus to let $r(\mathbf{z})$ as a standard Laplace distribution ($r(\mathbf{z}) = L(\mathbf{0}, I)$),

similarly, this yields the mean of L1 norm:

$$\mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log r(\mathbf{z})] \propto -\frac{1}{N} \sum_{i=1}^N \|z_i\|_1. \tag{17}$$

Based on Eq. 13, with Monte Carlo samples, our loss function is defined as:

$$\begin{aligned}
 \mathcal{L} &= -I_{NCE}(\mathbf{x}, \mathbf{z}) - \beta \cdot \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[\log r(\mathbf{z})] \\
 &\approx \frac{1}{N} \sum_{i=1}^N \left[-\log \frac{e^{f(x_i, z_i)}}{\frac{1}{K} \sum_{j=1}^K e^{f(x_i, z_j)}} + \beta \cdot \|z_i\|_p \right],
 \end{aligned} \tag{18}$$

where $p = 1$ or $p = 2$ according to the choice of r . We will explain the difference of this choice in the experiment part.

In this paper, we apply $f(x_i, z_i) = sim(E_\theta(x_i), E_\theta(\tilde{x}_i)) = sim(z_i, \tilde{z}_i)$, where z_i and \tilde{z}_i are the latent representations of x_i and its augmented view \tilde{x}_i ; $sim(u, v) \equiv u^\top \cdot v$ denotes the similarity between two normalized vectors u and v . The loss function can be reformulated as:

$$\mathcal{L} = \frac{1}{K} \sum_{i=1}^K \left[-\log \frac{\exp(sim(z_i, \tilde{z}_i))}{\frac{1}{K} \sum_{j=1}^K \exp(sim(z_i, \tilde{z}_j))} + \beta \|z_i\|_p \right], \tag{19}$$

Note that in practice the NCE lower bound requires larger number of negative samples to ensure the good performance [35], [41]. Considering the efficiency of negative sampling in anomaly detection, we follow SimCLR [35] and AMDIM [36] for the augmentation strategy. The pseudocode of the training process for the base model is shown in Algorithm 1.

Local Deep Infomax (local DIM) [32], which maximizes the mutual information between a global feature vector depend on the full input and a collection of local feature vectors pulled from an intermediate layer in the encoder, has shown to be greatly effective in improving feature learning and maximizing mutual information in [32], [36]. To get the best performance, we introduce local DIM to our method as the extension model. With the augmentation strategy from [35], [36], the negative samples size is increased and $L_{NCE}(z_i, \tilde{z}_i)$ becomes:

$$\mathcal{L}_{NCE}(z_i, \tilde{z}_i) = -\log \frac{\exp(sim(z_i, \tilde{z}_i))}{\frac{1}{2K} \sum_{j=1}^{2K} \mathbb{I}_{[j \neq i]} \exp(sim(z_i, \tilde{z}_j))}. \tag{20}$$

Thus the loss function, added local DIM, is formulated as:

$$\mathcal{L} = \frac{1}{K} \sum_{i=1}^K [\mathcal{L}_{NCE}(z_{g_i}, \tilde{z}_{g_i}) + \mathcal{L}_{NCE}(z_{g_i}, \tilde{z}_{l_i}) + \beta \cdot \|z_{g_i}\|_p], \tag{21}$$

where z_{g_i} refers to the global features, produced by parametric function E_θ , z_{l_i} refers to the local features, produced by an intermediate layer in E_θ .

D. Normal Score Measurement

Most surrogate supervision based approaches associate the anomaly score measurement with the loss function of surrogate task, which assume that anomalous data will result in relatively bigger loss. Following this mechanism, since the similarity

between \mathbf{z}_i and $\tilde{\mathbf{z}}_i$ is maximized in our method, a straight forward normal score measurement can be formulated as:

$$\text{NormalScore} = \text{sim}(z_i, \tilde{z}_i), \quad (22)$$

during z_i and \tilde{z}_i are the latent representations of the augmented view x_i and \tilde{x}_i of the same test data example x_{ori_i} , in which augmented view is randomly selected. In this design, sample with high normal score is consider to be normal. However, it is time consuming to traverse all combinations of augmented view pair and random selected one combination yield in unstable result. As a solution, when testing we skip the augmentation terms and encode the original data \mathbf{x}_{ori_i} directly. The corresponding output features is noted as \mathbf{z}_{ori_i} . Thus we reformulated the normal score for base model as:

$$\text{NormalScore} = \text{sim}(z_{ori_i}, z_{ori_i}). \quad (23)$$

The normal score for extension model is then reformulated as:

$$\text{NormalScore} = \text{sim}(z_{g_i}, z_{g_i}) + \text{sim}(z_{g_i}, z_{l_i}). \quad (24)$$

E. Relation to Other Algorithms

A successful theorem should not only be able to explain the implicit mechanisms in previous works but also be able to point out the potential improvement direction. We then take some classic methods to illustrate the interpretability of work.

AutoEncoder [42]: Reconstruction based anomaly detection using Auto-encoder is a mainstream method. It assumes that by minimizing the reconstruction error, normal and anomalous samples could lead to significantly different embedding and thus the corresponding reconstruction errors can be leveraged to differentiate the two types of samples. Then we will see how our equation fits this work. First we reformulate the mutual information terms as:

$$\begin{aligned} I_n(\mathbf{x}, \mathbf{z}) &= H_n(\mathbf{x}) - H_n(\mathbf{x}|\mathbf{z}) \\ &= H_n(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_n(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim p_n(\mathbf{z}|\mathbf{x})} [\log p_n(\mathbf{x}|\mathbf{z})]. \end{aligned} \quad (25)$$

With Eq. 25, the Eq. 10 is reformulated as:

$$\begin{aligned} &\max_{\theta} \{I_n(\mathbf{x}, \mathbf{z}) - H_n(\mathbf{z})\} \\ &= \max_{\theta} \{H_n(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_n(\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim p_n(\mathbf{z}|\mathbf{x})} [\log p_n(\mathbf{x}|\mathbf{z})] - H_n(\mathbf{z})\}, \end{aligned} \quad (26)$$

where the first term is a constant, the second term is the reconstruction likelihood, the third term is the entropy of \mathcal{Z} . This gives solid mathematical support to why anomaly detection can be benefited via maximizing reconstruction likelihood. Furthermore, it is indicated that adding entropy terms may further improve the method.

SVDD [4]: Deep one-class classification is the most classic method. Pre-trained by autoencoder network with reconstruction error, it then minimizes the volume of a data-enclosing hyper-sphere in latent space. It assumes that the normal examples of the data are closely mapped to center c , whereas anomalous examples are mapped further away from the center. SAD [31] mathematically proved that the objective function in SVDD minimizes an upper bound on the entropy of latent space. Considering the pre-trained autoencoding objective implicitly maximizes the mutual information, SAD intuitively summarizes

its objective function to have a positive correlation with mutual information and negative correlation with entropy,

which is consistent with our objective function in Eq. 10. our method gives theoretical support to why anomaly detection can benefit from this objective function. Both SAD and SVDD used a two-stage optimization method, which maximizes the mutual information first and then minimizes the entropy.

This two-stage optimization can not guarantee the simultaneous optimization of the mutual information and the entropy. Moreover, to relief the hyper-sphere collapse, only networks without bias terms and bounded activation functions can be used [4].

Information bottleneck [43], [44]: The Information Bottleneck (IB) methods define a good representation and learn it with the trade-off of a concise representation and a powerful prediction in downstream tasks [43], [44]. Various studies, extend the IB methods into deep learning scenario with variational methods, such as Variational Information Bottleneck (VIB) [45], Information Confidence Penalty (ICP) [46], Variational Discriminator Bottleneck (VDB) [47], *etc.* Our regularizer shown in Eq. 10 is consistent with the one proposed in [46]. We also show in Eq. 15 that the entropy is proportional to a KL divergence in terms of a reference distribution, which is consistent with the regularizer used in [45], [47].

Basic assumption for surrogate task based approaches in anomaly detection: As anomalous data is inaccessible during train process, most surrogate tasks based unsupervised anomaly detection methods are based on an assumption that data can be embedded into a lower-dimensional subspace where normal and anomalous samples appear significantly different [1]. In our method, we are the first to reveal the theoretical rationality for this basic assumption. From Proposition 2, we can see the key to reformulate the objective function from semi-supervised fashion (objective function in Eq. 4) to unsupervised fashion (objective function in Eq. 10) is to ensure the non-negativity of $\mathbb{E}_{p_n(\mathbf{x})} [H(p_n(\mathbf{z}|\mathbf{x}), p_a(\mathbf{z}|\mathbf{x}))]$. This is under the assumption, as introduced in Proposition 2, that $p_n(\mathbf{x}, \mathbf{z})$ and $p_a(\mathbf{x}, \mathbf{z})$ has certain distance such that for most samples $\mathbf{x}, \mathbf{z} \sim p_n(\mathbf{x}, \mathbf{z})$ the evaluated density $\log p_a(\mathbf{z}|\mathbf{x}) \leq 0$, which is consistent with the assumption in [1].

IV. EXPERIMENTS

In this section, we present a comprehensive set of experiments to validate our anomaly detection algorithm under unsupervised settings, in which multiple commonly used benchmark datasets are involved. To further demonstrate the robustness of our method, following the setting of ARNet [11], a subset of ImageNet [15] with higher resolution, richer texture and more complex background, is utilized.

A. Experimental Setups

Datasets. We experiment with the following five popular image datasets.

- MNIST [12]: consists of 70,000 28×28 handwritten grayscale digit images.
- Fashion-MNIST [13]: a relatively new dataset comprising 28×28 grayscale images of 70,000 fashion products from 10 categories, with 7,000 images per category.

- CIFAR-10 [14]: consists of 60,000 32×32 RGB images of 10 classes, with 6,000 images for per class. There are 50,000 training images and 10,000 test images, divided in a uniform proportion across all classes.
- CIFAR-100 [14]: consists of 100 classes, each of which contains 600 RGB images. The 100 classes in the CIFAR-100 are grouped into 20 “superclasses” to make the experiment more concise and data volume of each selected “normal class” larger.
- ImageNet [15]: Following ARNet [11], we group the data from the ILSVRC 2012 classification dataset [15] into 10 superclasses by merging similar category labels using Latent Dirichlet Allocation (LDA) [48], a natural language processing method (see appendix for more details). We noted this dataset is more challenging due to its higher resolution richer contexture and more complex background.

For all datasets, the training and test partitions remain as default. In addition, pixel values of all images are normalized to $[-1, 1]$.

Evaluation protocols. For a dataset with C classes, we conduct a batch of C experiments respectively with each of the C classes set as the “normal” class once. We then evaluate performance on an independent test set, which contains samples from all classes, including normal and anomalous data. As all classes have equal volumes of samples in our selected datasets, the overall number proportion of normal and anomalous samples is $1 : C - 1$. The model performance is then quantified using the area under the Receiver Operating Characteristic (ROC) curve metric (AUROC). It is commonly adopted as performance measurement in anomaly detection tasks and eliminates the subjective decision of threshold value to divide the “normal” samples from the anomalous ones. The above evaluation protocols are generally accepted among most recent works on anomaly detection [6]–[9], [11], [49]–[52].

Model configuration.

For all the datasets, the training epoch is set as 400, the learning rate is set as $2e-4$. The corresponding pseudocode of base model can be found in Algorithm 1.

For all experiments, if not specified, the base model refers to the model with loss function in Eq. 13, normal scoring function in Eq. 23, pseudocode in Algorithm 1. The extensive model refers to the model with loss function in Eq. 21, normal scoring function in Eq. 24. In general, we use a base model to investigate the property of the objective function and use an extensive model to reach better performance.

B. Investigation of Objective Function Properties

In this section, we look into the properties of our objective function with our base model. The objective function in Eq. 13 provides a lower bound for the maximization of the KL divergence in Eq. 2, which present a trade-off between the maximization of the mutual information $I(x, z)$ and the constrain of the entropy $H(z)$.

1) **Hyper parameter β in adjusting the trade-off:** To investigate the trade-off between the mutual information and the entropy, several experiments are conducted with a set of β , noted as $\mathcal{B} = \{0, 0.5, 1, 10, 20, 30, 40, 50, 60\}$, on datasets

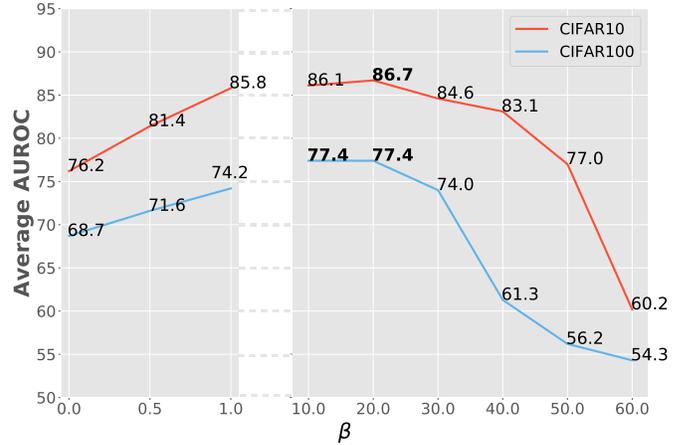


Fig. 1: Average AUROC w.r.t. β on CIFAR10 and CIFAR100.

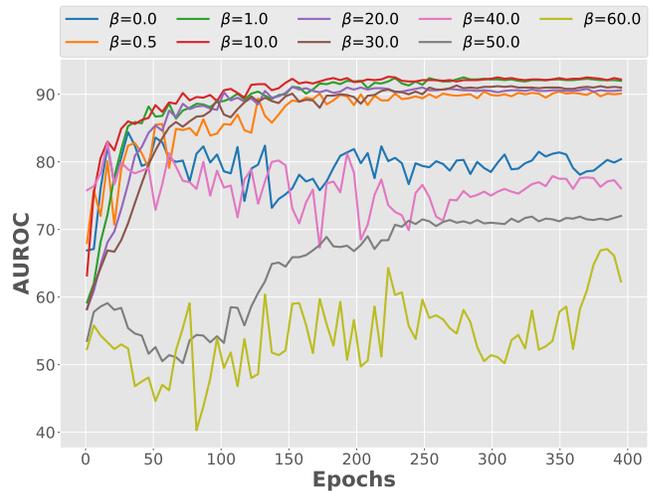


Fig. 2: The testing AUROC w.r.t. β during training.

CIFAR10 and CIFAR100. We report the average AUROC results of all classes with respect to β (Fig. 1). We observe that the model yield a comparable good performance with β in a relatively wide range from 0.5 to 40.0 in CIFAR10 and from 1 to 30.0 in CIFAR100, where the best performance is reached when $\beta = 20.0$ on both datasets, suggesting that the model is not very sensitive to the change in β as long as it in a certain range, which is vital in real-world cases where we can not utilize a testing dataset to select the super-parameter. In contrast, the model results in poor performance when β is either too small or too large.

In this work, if not specified, we choose $\beta = 20.0$ to conduct experiments on all datasets.

We also record the AUROC result every 5 epochs during the training process on CIFAR10 with different β in \mathcal{B} . The result of class 1 (car) is presented in Fig. 2. It is observed that the AUROC curve shows better convergence with β in the range from 0.5 to 30.0, where the model from the last epoch reaches almost the best performance. This is a vital property in real-world cases, where we can not utilize testing results to choose the model with the best performance.

To be noted, as observed from the above two experiments,

when $\beta = 0.0$ the model can not converge properly, which indicates that without the entropy constrain the unsupervised learning framework based on maximizing mutual information only is not directly applicable to anomaly detection.

2) **Digging out the essence of the trade-off:** We experiment on CIFAR10 and CIFAR100 with our base model under the hyper-parameter β in \mathcal{B} . For each training class, we investigate the relation between the converged training loss (mutual information and entropy) and the corresponding AUROC. Then we look into the learned representation z_i by calculating their distance from the center $\|z_i\|_2$ of normal and anomalous data as β increase. We report the results that trained with class 0 in Fig. 3. In the top panels of Fig. 3, we can observe that as β increases, the loss corresponds to the mutual information converges to a larger value while the loss corresponds to the entropy converges to a smaller value. As β increases, the model tends to ignore the mutual information, which matches our observation that the MI loss is getting larger and the entropy loss is getting smaller. It is remarkable that \mathcal{L}_{NCE} and \mathcal{L}_H become in similar scale when β is in the range from 0.5 to 40.0 in CIFAR10 and from 10.0 to 30.0 in CIFAR100, in which model results in better performance as shown in the top panels in Fig. 3. This indicates that a proper β for the trade-off should lead to similar converged scale for both \mathcal{L}_{NCE} and \mathcal{L}_H .

In the bottom panels of Fig. 3, we compare the mean $\|z_i\|_2$ of all samples between normal and anomalous data in testing datasets. As β increases, the mean $\|z_i\|_2$ for both normal and anomalous data shows a declining tendency. More importantly, the $\|z_i\|_2$ for the anomalous data decreases to a larger extent than that of the normal data. This results in a gap between normal and anomaly data. We observe that the model results in better performance when the mean $\|z_i\|_2$ gap is larger, with β in the range of 0.5 to 40.0 in CIFAR10 and 10.0 to 30.0 in CIFAR100. Especially, the biggest gap is attained when $\beta = 20.0$ on CIFAR10 and $\beta = 10.0$ on CIFAR100 dataset respectively, which correspond to their best performance as shown in the top panels of Fig. 3.

To analyze these experiments, we connect these results with the formulation. As shown in Eq. 15, the entropy regularizer is lower-bounded by the expected value of $\log r(\mathbf{z})$. As we choose r as a zero-centered reference distribution, whose density function is proportional to the p-norm $\|z_i\|_p$ ($p = 1$ or $p = 2$), a geometric interpretation is also guaranteed: the model will regularize the Euclidian (*resp.* Manhattan) distance from the center $\|z_i\|_2$ and encourage the representations z_i to be centered. From Fig. 3, as increase the value of β , we can observe that for anomaly data the mean $\|z_i\|_2$ is getting close to zero rapidly, while for normal data which profits from the mutual information maximization the mean $\|z_i\|_2$ is getting close to zero relatively slower. When β is too large (*e.g.* $\beta \geq 60.0$), we can observe that the maximization of MI is over-regularized, thus the model cannot extract effective features for normal data, which results in a lower AUROC. These results enlighten the importance of the trade-off between mutual information and entropy.

C. Ablation: Mutual Information and Entropy Estimators

The previous section discusses the importance of the information trade-off. A nice property of our formulation is that Eq. 13

is general enough and is able to plug in alternative estimators for the mutual information and entropy loss. In this section, we conduct the ablation experiments with the alternative mutual information and entropy estimator and justify the expressiveness of the estimators that we choose.

1) **Learning with a different mutual information estimator:** We apply the InfoNCE estimator [40] for the calculation of the mutual information. Considering there are a wide range of approaches for the MI estimation, we also consider an alternative lower bound for MI estimation called Donsker-Varadahn (DV) bound [53].

$$I(\mathbf{x}; \mathbf{z}) \geq \mathbb{E}_{p(\mathbf{x}, \mathbf{z})}[f(\mathbf{x}, \mathbf{z})] - \log \mathbb{E}_{p(\mathbf{z})}[e^{f(\mathbf{x}, \mathbf{z})}] \triangleq I_{DV} \quad (27)$$

This bound is widely used in various algorithms such as Mutual Information Neural Estimator (MINE) [33]. Deep InfoMax (DIM) [32] continues to extend the critic in this bound with a Jensen-Shannon based formulation as:

$$I(\mathbf{x}; \mathbf{z}) \geq \mathbb{E}[-\sigma(-f(\mathbf{x}, \mathbf{z})) - \sigma(f(\mathbf{x}', \mathbf{z}))] \triangleq I_{JSD} \quad (28)$$

where $\sigma(z) = \log(1 + e^z)$ denotes the Softplus function, and note that for DV bound and JSD bound the critic function f is usually modeled with a discriminator function $f: \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$. Moreover, according to existing studies like [32], the JSD-based estimator and the DV-based estimator behave similarly. Here, we apply the JSD-based MI estimator to represent the effect of DV bound.

To compare the JSD-based and NCE-based MI estimator, on our base model, we utilize the JSD-based MI estimator and the InforNCE estimator. The experiments are conducted on CIFAR10 and CIFAR100 and the hyper-parameter is set as $\beta = 20.0$. The results are shown in Table I. Given the same entropy estimator, we can find the model with the JSD-based MI estimator underperforms than the NCE-based MI estimator.

To investigate the performance gap, we investigate the training process and illustrate the testing AUROC curve and the mutual information loss during the training process in Fig. 4, where the NCE-based and the JSD-based estimator is marked in red and blue respectively. Compared to the NCE-based estimator, both testing AUROC and the loss by JSD-based estimator show large perturbation during the training, which makes the model more difficult to converge.

One plausible explanation is that the mutual information estimation variance of the JSD estimator is much larger than the NCE estimator, which is the critical difference between these two MI estimators [41]. This large variance caused by the JSD-based estimator makes the model more unstable during the training and the learned representations may be useless. In this way, the JSD-based methods need more fine-tuning to stabilize the training process of the model.

2) **Learning with a different entropy estimator:** We introduce two reference distributions in Eq. 16 and Eq. 17 to estimate the entropy, corresponding to the L2 norm and the L1 norm, respectively. Similar to the previous experiments, CIFAR10 and CIFAR100 are used. With different entropy estimators, we test with both the NCE-based and the JSD-based mutual information estimator, and report the mean and standard deviation of the AUROC among all classes in Table I. Given the same MI estimator, the L1 norm estimator yields

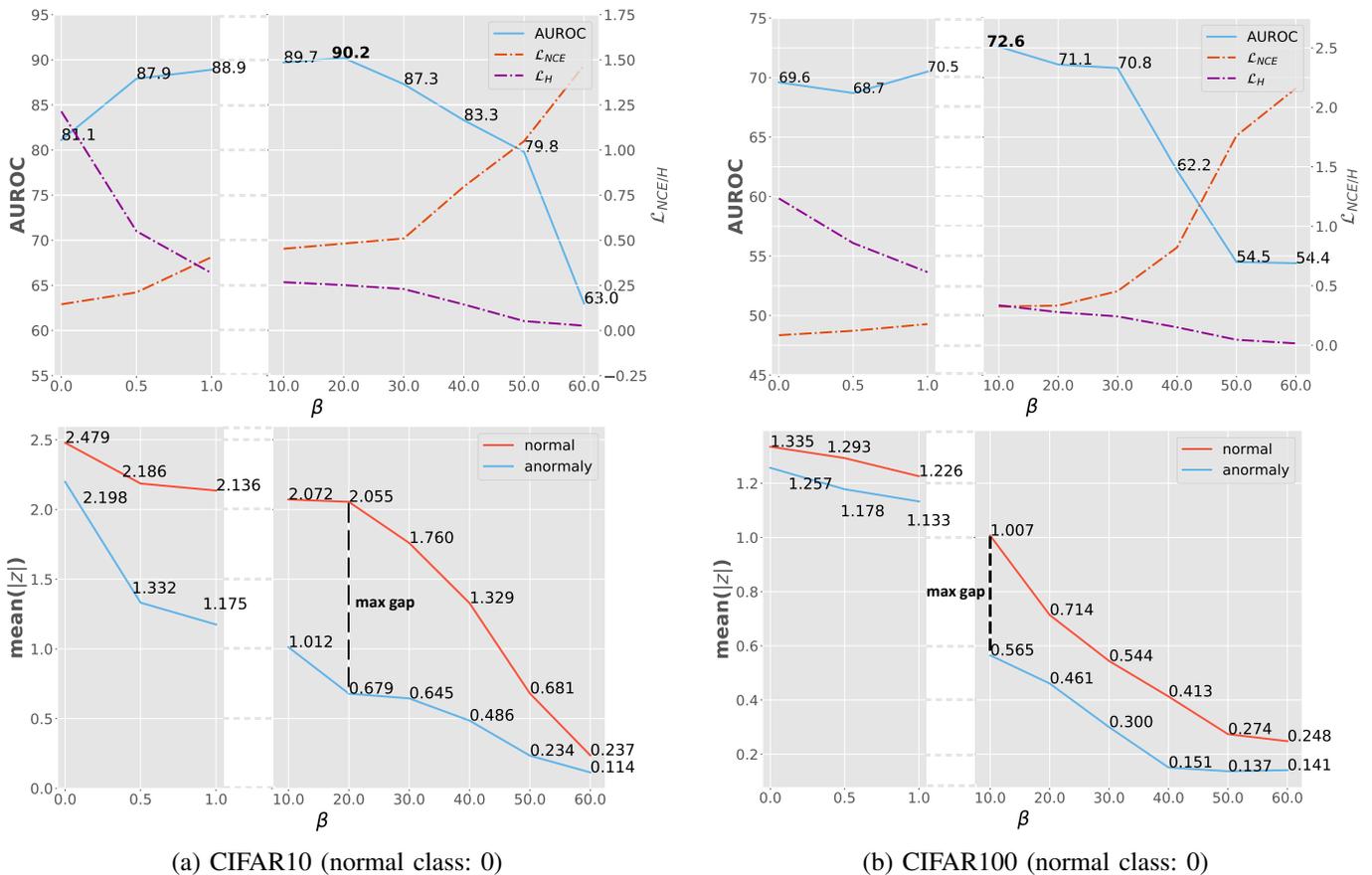


Fig. 3: The trade-off between MI and entropy. (Top) The AUROC curve (solid line in blue), the loss curve of NCE (dashed dot line in red) and the loss curve of entropy (dashed dot line in purple) w.r.t. β . (Bottom) The mean value of $\|z\|$ w.r.t. β that corresponds to normal and anomalous data. (Best viewed in color)

TABLE I: Comparison of different MI and entropy estimators.

Loss function	CIFAR10		CIFAR100	
	AVG	STD	AVG	STD
$\mathcal{L}_{NCE} + \beta \cdot \mathcal{L}_2 Entropy$	85.5	7.02	76.9	5.26
$\mathcal{L}_{NCE} + \beta \cdot \mathcal{L}_1 Entropy$	86.7	5.47	77.4	4.62
$\mathcal{L}_{JSD} + \beta \cdot \mathcal{L}_2 Entropy$	72.5	7.20	69.6	4.73
$\mathcal{L}_{JSD} + \beta \cdot \mathcal{L}_1 Entropy$	73.9	5.00	72.7	3.21

TABLE II: Efficient normal scoring mechanism.

Normal scoring function	CIFAR10		CIFAR100	
	AVG	STD	AVG	STD
$NormalScore_{rand}$	86.4	1.04	75.6	1.26
$NormalScore_{mc}$	86.9	0.05	78.5	0.12
$NormalScore_{ori}$	86.7	—	77.4	—

better performance. A case-by-case study is conducted and based on the hyper-parameter that we set, when the L1 norm used as the entropy regularizer, the representations are more inclined to shrink to the center. This results in a larger distance

gap between normal and anomaly data than the L2 norm does and in order to have the same performance L2 norm may need more proper hyper-parameter tuning to balance the information trade-off.

Based on our observation, the NCE-based estimator is more suitable and compatible with the proposed entropy estimator. In the following sections, we continue to refine the NCE-based estimator and apply the L1 entropy regularizer to present the effectiveness of our proposed method.

D. Efficient Normal Scoring Mechanism

As shown in Eq. 22 the normal score is firstly designed to be the similarity of z_i and \tilde{z}_i , which are the latent representations of the randomly augmented view x_i and \tilde{x}_i of the same test data example x_{ori} . However, the random augmentation yield in unstable result. An alternative solution is to utilized Monte Carlo samples, and the normal score is reformulated as:

$$NormalScore = \sum_{h=1}^H [sim(z_{g_i}^{(h)}, \tilde{z}_{g_i}^{(h)})] \quad (29)$$

Where the h refers to the h^{st} sampling. Despite being more stable, it is time-consuming. We introduce a compromise as shown in Eq. 23, which skip the augmentation and directly input the original sample.

TABLE III: Average area under the ROC curve (AUROC) in % of anomaly detection methods. For every dataset, each model is trained on the single class, and tested against all other classes. ‘‘SD’’ means standard deviation among classes. The best performing method is in bold. Results of Deep SVDD are borrowed from [4]

. Results of VAE, AnoGAN and ADGAN are from [6]. Results of DAGMM, DSEBM and GeoTrans are from [9]. Results of GANomaly and ARNet are from [11].

Dataset	Method	0	1	2	3	4	5	6	7	8	9	avg	std
MNIST	VAE [49]	92.1	99.9	81.5	81.4	87.9	81.1	94.3	88.6	78.0	92.0	87.7	7.05
	D-SVDD [4]	98.0	99.7	91.7	91.9	94.9	88.5	98.3	94.6	93.9	96.5	94.8	3.46
	AnoGAN [52]	99.0	99.8	88.8	91.3	94.4	91.2	92.5	96.4	88.3	95.8	93.7	4.00
	ADGAN [6]	99.5	99.9	93.6	92.1	94.9	93.6	96.7	96.8	85.4	95.7	94.7	4.15
	GANomaly [7]	97.2	99.6	85.1	90.6	94.9	94.9	97.1	93.9	79.7	95.4	92.8	6.12
	OCGAN [8]	99.8	99.9	94.2	96.3	97.5	98.0	99.1	98.1	93.9	98.1	97.5	2.10
	GeoTrans [9]	98.2	91.6	99.4	99.0	99.1	99.6	99.9	96.3	97.2	99.2	98.0	2.50
	ARNet [11]	98.6	99.9	99.0	99.1	98.1	98.1	99.7	99.0	93.6	97.8	98.3	1.78
	Our Base model	91.3	98.6	90.6	88.6	88.2	89.1	91.2	85.6	87.7	86.0	89.7	3.70
Our Extension model	99.5	99.7	98.8	98.3	97.7	96.7	98.7	97.5	98.6	97.3	98.3	0.97	
Fashion-MNIST	DAGMM [50]	42.1	55.1	50.4	57.0	26.9	70.5	48.3	83.5	49.9	34.0	51.8	16.47
	DSEBM [51]	91.6	71.8	88.3	87.3	85.2	87.1	73.4	98.1	86.0	97.1	86.6	8.61
	ADGAN [6]	89.9	81.9	87.6	91.2	86.5	89.6	74.3	97.2	89.0	97.1	88.4	6.75
	GANomaly [7]	80.3	83.0	75.9	87.2	71.4	92.7	81.0	88.3	69.3	80.3	80.9	7.37
	GeoTrans [9]	99.4	97.6	91.1	89.9	92.1	93.4	83.3	98.9	90.8	99.2	93.5	5.22
	ARNet [11]	92.7	99.3	89.1	93.6	90.8	93.1	85.0	98.4	97.8	98.4	93.9	4.70
	Our Base model	94.3	99.1	91.6	95.3	90.9	98.0	86.1	98.0	97.2	97.1	94.8	4.11
	Our Extension model	96.7	99.7	95.3	97.3	95.1	99.2	89.8	99.3	99.1	99.3	97.1	3.08
	CIFAR-10	VAE [49]	62.0	66.4	38.2	58.6	38.6	58.6	56.5	62.2	66.3	73.7	58.1
D-SVDD [4]		61.7	65.9	50.8	59.1	60.9	65.7	67.7	67.3	75.9	73.1	64.8	7.16
DAGMM [50]		41.4	57.1	53.8	51.2	52.2	49.3	64.9	55.3	51.9	54.2	53.1	5.95
DSEBM [51]		56.0	48.3	61.9	50.1	73.3	60.5	68.4	53.3	73.9	63.6	60.9	9.10
AnoGAN [52]		61.0	56.5	64.8	52.8	67.0	59.2	62.5	57.6	72.3	58.2	61.2	5.68
ADGAN [6]		63.2	52.9	58.0	60.6	60.7	65.9	61.1	63.0	74.4	64.4	62.4	5.56
GANomaly [7]		93.5	60.8	59.1	58.2	72.4	62.2	88.6	56.0	76.0	68.1	69.5	13.08
OCGAN [8]		75.7	53.1	64.0	62.0	72.3	62.0	72.3	57.5	82.0	55.4	65.6	9.52
GeoTrans [9]		74.7	95.7	78.1	72.4	87.8	87.8	83.4	95.5	93.3	91.3	86.0	8.52
ARNet [11]	78.5	89.8	86.1	77.4	90.5	84.5	89.2	92.9	92.0	85.5	86.6	5.35	
ImageNet	Our Base model	89.9	95.7	80.6	72.8	80.9	83.8	91.1	91.1	93.0	87.9	86.7	5.47
	Our Extension model	93.0	97.6	88.5	85.8	89.5	92.1	95.2	93.6	95.0	95.5	92.6	3.65
	GANomaly [7]	58.9	57.5	55.7	57.9	47.9	61.2	56.8	58.2	49.7	48.8	55.3	4.46
ImageNet	GeoTrans [9]	72.9	61.0	66.8	82.0	56.7	70.1	68.5	77.2	62.8	83.6	70.1	8.43
	ARNet [11]	71.9	85.8	70.7	78.8	69.5	83.3	80.6	72.4	74.9	84.3	77.2	5.77
	Our Base model	79.8	81.3	80.7	85.4	79.7	84.1	75.3	81.5	76.1	67.1	79.1	5.23
Our Extension model	88.2	89.8	86.1	89.6	89.4	89.4	80.3	85.7	81.8	74.9	85.5	5.02	
Dataset	Method	0	1	2	3	4	5	6	7	8	9	10	
CIFAR-100	DAGMM [50]	43.4	49.5	66.1	52.6	56.9	52.4	55.0	52.8	53.2	42.5	52.7	
	DSEBM [51]	64.0	47.9	53.7	48.4	59.7	46.6	51.7	54.8	66.7	71.2	78.3	
	ADGAN [6]	63.1	54.9	41.3	50.0	40.6	42.8	51.1	55.4	59.2	62.7	79.8	
	GANomaly [7]	57.9	51.9	36.0	46.5	46.6	42.9	53.7	59.4	63.7	68.0	75.6	
	GeoTrans [9]	74.7	68.5	74.0	81.0	78.4	59.1	81.8	65.0	85.5	90.6	87.6	
	ARNet [11]	77.5	70.0	62.4	76.2	77.7	64.0	86.9	65.6	82.7	90.2	85.9	
	Our Base model	71.1	80.2	78.8	79.5	78.9	79.6	78.5	75.6	74.3	84.8	85.2	
	Our Extension model	82.0	85.8	87.8	86.2	89.3	87.6	86.8	83.0	85.1	91.3	90.1	
	Method	11	12	13	14	15	16	17	18	19	avg	SD	
DAGMM [50]	46.4	42.7	45.4	57.2	48.8	54.4	36.4	52.4	50.3	50.5	6.55		
DSEBM [51]	62.7	66.8	52.6	44.0	56.8	63.1	73.0	57.7	55.5	58.8	9.36		
ADGAN [6]	53.7	58.9	57.4	39.4	55.6	63.3	66.7	44.3	53.0	54.7	10.08		
GANomaly [7]	57.6	58.7	59.9	43.9	59.9	64.4	71.8	54.9	56.8	56.5	9.94		
GeoTrans [9]	83.9	83.2	58.0	92.1	68.3	73.5	93.8	90.7	85.0	78.7	10.76		
ARNet [11]	83.5	84.6	67.6	84.2	74.1	80.3	91.0	85.3	85.4	78.8	8.82		
Our Base model	73.4	74.2	72.1	72.8	73.1	69.8	84.0	80.3	81.0	77.4	4.62		
Our Extension model	80.9	84.2	80.3	85.2	79.7	82.2	93.9	91.5	87.9	86.0	3.98		

With the base model, we conduct experiment on CIFAR10 and CIFAR100 to evaluate the three normal scoring mechanism based on Eq. 22, 23 and 29, which is noted as $NormalScore_{rand}$, $NormalScore_{ori}$ and $NormalScore_{mc}$

accordingly. To be noted that, in Eq. 29, H is set as 100. We repeat the testing process for 10 times for each normal scoring mechanism and record the average and standard deviation of the AUROC of all classes in CIFAR10. As illustrated in

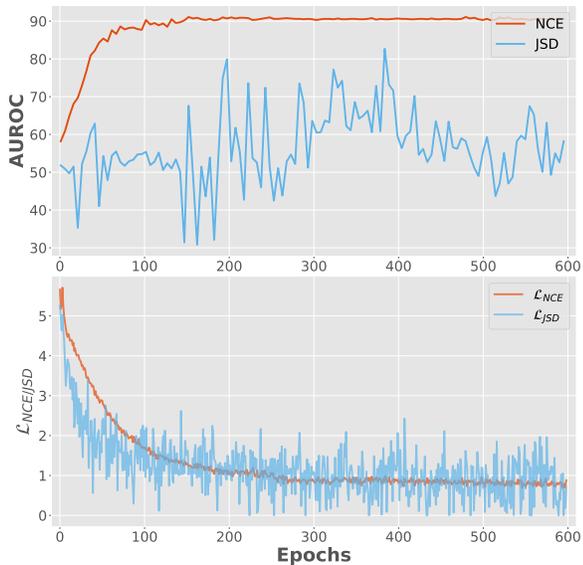


Fig. 4: The comparison of NCE-based and JSD-based mutual information estimator: conducted on dataset CIFAR10 where class 1 (car) is considered as normal data. **(Top)** The testing ROC curve during the training process by using NCE and JSD estimator for mutual information calculation. **(Bottom)** The estimated loss by using the NCE-based and JSD-based mutual information estimator during the training process.

Table II, the normal scoring mechanisms base on Eq. 23 yield a comparable and stable result with only one inference. We take Eq. 23/Eq. 24 as our refined normal scoring mechanism for the base/extension model through in the following experiments to obtain the best performance.

E. Extensive Experiment on Extension Model

Table III provides results of our base model and extension model on MNIST, Fashion-MNIST, CIFAR-10, ImageNet and CIFAR-100. For grayscale datasets, such as MNIST and Fashion-MNIST, we convert the image into 3 channels through expanding. The performance of extension model is improved considerably over the base model both in average and standard deviation of the AUROC with Local DIM [36], indicating that a more effective module to optimize mutual information is beneficial. On all involved datasets, experiment results illustrate that the average AUROC or the standard deviation of our method outperforms all other methods to different extents. Furthermore, our method reveals greater advantage in more difficult datasets, from FashionMNIST which surpasses the SOTA by 3.2%, to 6% in CIFAR10, to 7.2% in CIFAR100 and to 8.3% in ImageNet Subset. This may be because our method is benefit from a decoder-free representation learning framework, where an advanced feature extracting model like RESNET can be utilized to handle pictures with higher resolution and more complex texture. More importantly, our method results in the lower standard deviation of AUROCs, which reflect the stability of the model when dealing with different kinds of anomalous data. This is vital especially in anomaly detection, where anomalous data can not be foreseen [11].

V. CONCLUSION AND FUTURE WORK

In this paper, we are the first to put forward an anomaly detection based objective function that can be optimized end-to-end in an unsupervised fashion. Many works can find mathematical support and further discover the potential optimization direction through our method. At last, based on the objective function we present a method that overperforms the state-of-the-arts, which illustrates the correctness of our objective function and rationality of the design of the loss function. Despite we restrict our method in an unsupervised anomaly detection setting, as can be seen in Equation 4, this work can be extended to semi-supervised anomaly detection, which remains to be our future work. Notably, there are other loss functions in maximizing mutual information and minimizing entropy to explore. Furthermore, we can investigate the specific functions of each component in our lower bound objective function in Eq. 13. The function of mutual information seems obvious, as widely explored in unsupervised representation learning approaches [32], [34]–[37], maximizing mutual information can effectively force the model to obtain better representation. Since the entropy forces $\|z_i\|_2$ to be closed to zero, it may be indicated that the entropy regularizer limits the model representation ability. Thus this trade-off seems to force the network to represent normal data with limited representation ability. In another word, this trade-off force the network to generate features only to properly represent normal data, which is consistent with the insight in [8]. As anomalous data is inaccessible during training, one feasible solution to enlarge the distribution of normal and anomalous data in latent space is to make anomalous data unable to be represented properly. Is this the essence of anomaly detection? We will explore this in our future work. .

REFERENCES

- [1] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly detection: A survey,” *Acm Computing Surveys*, vol. 41, no. 3, 2009.
- [2] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, “Estimating the support of a high-dimensional distribution,” *Neural computation*, vol. 13, no. 7, pp. 1443–1471, 2001.
- [3] P. Oza and V. M. Patel, “One-class convolutional neural network,” *IEEE Signal Processing Letters*, vol. 26, no. 2, pp. 277–281, 2018.
- [4] L. Ruff, N. Görnitz, L. Deecke, S. A. Siddiqui, R. Vandermeulen, A. Binder, E. Müller, and M. Kloft, “Deep one-class classification,” in *International Conference on Machine Learning*, 2018.
- [5] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli, “Adversarially learned one-class classifier for novelty detection,” in *The IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [6] L. Deecke, R. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft, “Anomaly detection with generative adversarial networks,” *arXiv preprint arXiv:1809.04758*, 2018.
- [7] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon, “GANomaly: Semi-Supervised Anomaly Detection via Adversarial Training,” *Asian Conference on Computer Vision*, 2018.
- [8] P. Perera, R. Nallapati, and B. Xiang, “Ocgan: One-class novelty detection using gans with constrained latent representations,” in *The IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [9] I. Golan and R. El-Yaniv, “Deep anomaly detection using geometric transformations,” in *Advances in Neural Information Processing Systems*, 2018.
- [10] L. X. e. a. Wang S, Zeng Y, “Effective end-to-end unsupervised outlier detection via inlier priority of discriminative network,” in *Advances in Neural Information Processing Systems*, 2019, pp. 5962–5975.
- [11] F. Ye, C. Huang, J. Cao, M. Li, Y. Zhang, and C. Lu, “Attribute restoration framework for anomaly detection,” *arXiv preprint arXiv:1911.10676v2*, 2020.

- [12] Y. LeCun, "The mnist database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>, 1998.
- [13] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [14] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [15] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, 2015.
- [16] C. S. Chalapathy R, "Deep learning for anomaly detection: A survey," *arXiv preprint arXiv:1901.03407*, 2019., 2019.
- [17] S. S. Markou M, "Novelty detection: a review-part 1: statistical approaches," *Signal Processing*, 2003.
- [18] —, "Novelty detection: a review-part 2: neural network based approaches," *Signal Processing*, 2003.
- [19] M. A. F. Pimentel, D. A. Clifton, C. Lei, and L. Tarassenko, "A review of novelty detection," *Signal Processing*, 2014.
- [20] B. R. Kiran, D. M. Thomas, and R. Parakkal, "An overview of deep learning based methods for unsupervised and semi-supervised anomaly detection in videos," *Journal of Imaging*, vol. 4, no. 2, p. 36, 2018.
- [21] W. Chu, H. Xue, C. Yao, and D. Cai, "Sparse coding guided spatiotemporal feature learning for abnormal event detection in large videos," *IEEE Transactions on Multimedia*, vol. 21, no. 1, pp. 246–255, 2018.
- [22] K. Xu, X. Jiang, and T. Sun, "Anomaly detection based on stacked sparse coding with intraframe classification strategy," *IEEE Transactions on Multimedia*, vol. 20, no. 5, pp. 1062–1074, 2018.
- [23] K. Xu, T. Sun, and X. Jiang, "Video anomaly detection and localization based on an adaptive intra-frame classification network," *IEEE Transactions on Multimedia*, 2019.
- [24] E. Eskin, "Anomaly detection over noisy data using learned probability distributions," in *International Conference on Machine Learning*, 2000.
- [25] K. Yamanishi, J. I. Takeuchi, G. Williams, and P. Milne, "On-line unsupervised outlier detection using finite mixtures with discounting learning algorithms," *Data Mining & Knowledge Discovery*, 2000.
- [26] M. Rahmani and G. K. Atia, "Coherence pursuit: Fast, simple, and robust principal component analysis," *IEEE Transactions on Signal Processing*, vol. 65, no. 23, pp. 6260–6275, 2017.
- [27] H. Xu, C. Caramanis, and S. Sanghavi, "Robust pca via outlier pursuit," *IEEE Transactions on Information Theory*, vol. 58, no. 5, pp. 3047–3064, 2012.
- [28] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, "Autoencoding beyond pixels using a learned similarity metric," *PMLR*, pp. 1558–1566, 2016.
- [29] A. Dosovitskiy and T. Brox, "Generating images with perceptual similarity metrics based on deep networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 658–666.
- [30] D. Hendrycks, M. Mazeika, and T. Dietterich, "Deep anomaly detection with outlier exposure," in *International Conference on Learning Representations*, 2019.
- [31] L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K.-R. Müller, and M. Kloft, "Deep semi-supervised anomaly detection," *International Conference on Learning Representations*, 2020.
- [32] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," in *International Conference on Learning Representations*, 2019.
- [33] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm, "Mutual information neural estimation," in *International Conference on Machine Learning*, 2018, pp. 531–540.
- [34] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [35] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," *arXiv preprint arXiv:2002.05709*, 2020.
- [36] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," in *Advances in Neural Information Processing Systems*, 2019, pp. 15 535–15 545.
- [37] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.
- [38] T. M. Cover and J. A. Thomas, "Elements of information theory 2nd edition (Wiley series in Telecommunications and Signal Processing)," 2006.
- [39] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville, "Adversarially learned inference," in *International Conference on Learning Representations*, 2017.
- [40] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010, pp. 297–304.
- [41] B. Poole, S. Ozair, A. van den Oord, A. Alemi, and G. Tucker, "On variational bounds of mutual information," in *International Conference on Machine Learning*, 2019.
- [42] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *Proceedings of the 25th international conference on Machine learning*, 2008, pp. 1096–1103.
- [43] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," *arXiv preprint physics/0004057*, 2000.
- [44] N. Tishby and N. Zaslavsky, "Deep learning and the information bottleneck principle," in *2015 IEEE Information Theory Workshop (ITW)*. IEEE, 2015, pp. 1–5.
- [45] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," in *International Conference on Learning Representations*, 2017.
- [46] G. Pereyra, G. Tucker, J. Chorowski, Ł. Kaiser, and G. Hinton, "Regularizing neural networks by penalizing confident output distributions," in *Workshop Track in International Conference on Learning Representations*, 2017.
- [47] X. B. Peng, A. Kanazawa, S. Toyer, P. Abbeel, and S. Levine, "Variational discriminator bottleneck: Improving imitation learning, inverse rl, and gans by constraining information flow," in *International Conference on Learning Representations*, 2018.
- [48] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, pp. 933–1022, 2003.
- [49] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [50] S. Zhai, Y. Cheng, W. Lu, and Z. Zhang, "Deep structured energy based models for anomaly detection," in *International Conference on Machine Learning*, 2016.
- [51] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen, "Deep autoencoding gaussian mixture model for unsupervised anomaly detection," in *International Conference on Learning Representations*, 2018.
- [52] 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 *International Conference on Information Processing in Medical Imaging*, 2017.
- [53] M. Donsker and S. Varadhan, "Large deviations for markov processes and the asymptotic evaluation of certain markov process expectations for large times," in *Probabilistic Methods in Differential Equations*. Springer, 1975, pp. 82–88.

APPENDIX A
MODEL STRUCTURE

The detailed structure of the model we used can be found in this section, where the model structure in Table V is utilized for dataset ImageNet and model structure in Table IV is utilized for other datasets.

TABLE IV: Small Encoder Architecture

Small Encoder Architecture	
ReLU(Conv 2d(3, ndf, 3, 1, 0))	
ResBlock (ndf, ndf, 1, 1, 0)	
ResBlock (1*ndf, 2*ndf, 4, 2, ndepth)	} – provides f_{local}
ResBlock (2*ndf, 4*ndf, 2, 2, ndepth)	
ResBlock (4*ndf, 4*ndf, 3, 1, ndepth)	} – provides f_{global}
ResBlock (4*ndf, 4*ndf, 3, 1, ndepth)	
ResBlock (4*ndf, nrkhs , 3, 1, 1)	
$ndf = 128, nrkhs = 1024, ndepth = 10$	

TABLE V: Big Encoder Architecture

Big Encoder Architecture	
ReLU(Conv 2d(3, ndf, 5, 2, 2))	
ReLU(Conv 2d(ndf, ndf, 3, 1, 0))	
ResBlock (1*ndf, 2*ndf, 4, 2, ndepth)	} – provides f_{local}
ResBlock (2*ndf, 4*ndf, 4, 2, ndepth)	
ResBlock (4*ndf, 8*ndf, 2, 2, ndepth)	} – provides f_{global}
ResBlock (8*ndf, 8*ndf, 3, 1, ndepth)	
ResBlock (8*ndf, 8*ndf, 3, 1, ndepth)	
ResBlock (8*ndf, nrkhs , 3, 1, 1)	
$ndf = 192, nrkhs = 1536, ndepth = 8$	

APPENDIX B

PSEUDOCODE FOR EXTENSION METHOD

The testing pseudocode of base model can be found in Algorithm 2. The training and testing pseudocode of extension model can be found in Algorithm 4, 5. The NCE Loss pseudocode of extension model can be found in Algorithm 3.

Algorithm 2: Testing Pseudocode for base model

```

Input: batch size N, similarity function  $sim$ , structure of model  $E_\theta$ 
1 for sampled minibatch  $\{x_k\}_{k=1}^N$  do
2   for all  $k \in \{1, \dots, N\}$  do
3     # Extract features
4      $z_{ki} = E_\theta(x_i)$ 
5   for all  $i \in \{1, \dots, N\}$  do
6     # Similarity
7      $s_i = sim(z_i, z_i) = z_i^\top \cdot z_i$ 
8      $s'_i = c_2 \cdot \tanh\left(\frac{s_i}{c_1 \cdot c_2}\right)$ 
9     #normal data has larger  $s'_i$ 
10 return  $roc\_auc\_score(s'_i, label_i)$ 

```

Algorithm 3: NCE Loss Pseudocode extension model

```

Input: batch size N, input feature  $\Phi_1, \Phi_2$ 
1 # shape  $\Phi_1(n\_batch, n\_dim, n_1)$ 
2 # shape  $\Phi_2(n\_batch, n\_dim, n_2)$ 
3 for sampled minibatch  $\{x_k\}_{k=1}^N$  do
4   for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
5     # shape  $\Phi_1(n\_dim, n_1)$ 
6     # shape  $\Phi_2(n\_dim, n_2)$ 
7     # Similarity
8      $s_{i,j} = sim(\phi_{1i}, \phi_{2j}) = \sum_{n_1} \sum_{n_2} \phi_{1i}^\top \cdot \phi_{2j}$ 
9      $s'_{i,j} = c_2 \cdot \tanh\left(\frac{s_{i,j}}{c_1 \cdot c_2}\right)$ 
10  define  $S = \{s'_{i,j}, i, j \in 2N\}$ ,  $s_{max} = \max\{S\}$ 
11  define  $S_{shift} = \{\hat{s}_{i,j} = s'_{i,j} - s_{max}, i, j \in 2N\}$ 
12  define  $\ell(i, j) = -\log \frac{\exp(\hat{s}_{i,j})}{\sum_{k=1}^{2N} \mathbb{I}_{[k \neq i]} \exp(\hat{s}_{i,k})}$ 
13   $\mathcal{L}_{nce} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
14 return  $\mathcal{L}_{nce}$ 

```

Algorithm 4: Training Pseudocode extension model

```

Input: batch size N, NCE Loss function  $\mathcal{L}_{nce}$ , encoder model  $E_{\theta_1}$ ,
nonlinear function  $\phi_{\theta_2}$ , entropy weight  $\beta$ 
1 for sampled minibatch  $\{x_k\}_{k=1}^N$  do
2   for all  $k \in \{1, \dots, N\}$  do
3     randomly draw two augmentation functions  $t, t' \sim \mathcal{T}$ 
4     # Augmentation
5      $\tilde{x}_{2k-1} = t(x_k)$ 
6      $\tilde{x}_{2k} = t'(x_k)$ 
7     # Extract features
8      $g_{2k-1}, l_{2k-1} = E_{\theta_1}(\tilde{x}_{2k-1})$ 
9      $g_{2k}, l_{2k} = E_{\theta_1}(\tilde{x}_{2k})$ 
10    # nonlinear projection
11     $\tilde{l}_{2k-1} = \phi_{\theta_2}(l_{2k-1})$ 
12     $l_{2k} = \phi_{\theta_2}(l_{2k})$ 
13     $\mathcal{L}_{gvg} = \mathcal{L}_{nce}(g_{2k-1}, g_{2k})$ 
14     $\mathcal{L}_{gvl} = 0.5 * [\mathcal{L}_{nce}(g_{2k}, l_{2k-1})] + \mathcal{L}_{nce}(g_{2k-1}, l_{2k})$ 
15     $\mathcal{L}_{entropy} = \frac{1}{2N} \sum_{k=1}^{2N} (\|g_i\|_q + \|l_i\|_q)$ 
16     $\mathcal{L} = \mathcal{L}_{gvg} + \mathcal{L}_{gvl} + \beta \cdot \mathcal{L}_{entropy}$ 
17    update networks  $E_\theta$ 
18 return encoder network  $E_\theta$ 

```

Algorithm 5: Testing Pseudocode extension model

```

Input: batch size N, similarity function  $sim$ , encoder model  $E_{\theta_1}$ ,
nonlinear function  $\phi_{\theta_2}$ 
1 for sampled minibatch  $\{x_k\}_{k=1}^N$  do
2   for all  $k \in \{1, \dots, N\}$  do
3     # Extract features
4      $g_k, l_k = E_{\theta_1}(x_k)$ 
5     # nonlinear projection
6      $\tilde{l}_k = \phi_{\theta_2}(l_k)$ 
7   for all  $k \in \{1, \dots, N\}$  do
8     # Similarity
9      $s_{gvg_k} = sim(g_k, g_k)$ 
10     $s_{gvl_k} = sim(g_k, l_k)$ 
11     $s_k = s_{gvg_k} + s_{gvl_k}$ 
12     $s'_k = c_2 \cdot \tanh\left(\frac{s_k}{c_1 \cdot c_2}\right)$ 
13    #normal data has larger  $s'_i$ 
14 return  $roc\_auc\_score(s'_i, label_i)$ 

```