

Unsupervised Anomaly Detection for Tabular Data Using Noise Evaluation

Wei Dai, Kai Hwang, Jicong Fan*

The Chinese University of Hong Kong, Shenzhen, China
weidai@link.cuhk.edu.com, hwangkai@cuhk.edu.cn, fanjicong@cuhk.edu.cn

Abstract

Unsupervised anomaly detection (UAD) plays an important role in modern data analytics and it is crucial to provide simple yet effective and guaranteed UAD algorithms for real applications. In this paper, we present a novel UAD method for tabular data by evaluating how much noise is in the data. Specifically, we propose to learn a deep neural network from the clean (normal) training dataset and a noisy dataset, where the latter is generated by adding highly diverse noises to the clean data. The neural network can learn a reliable decision boundary between normal data and anomalous data when the diversity of the generated noisy data is sufficiently high so that the hard abnormal samples lie in the noisy region. Importantly, we provide theoretical guarantees, proving that the proposed method can detect anomalous data successfully, although the method does not utilize any real anomalous data in the training stage. Extensive experiments through more than 60 benchmark datasets demonstrate the effectiveness of the proposed method in comparison to 12 baselines of UAD. Our method obtains a 92.27% AUC score and a 1.68 ranking score on average. Moreover, compared to the state-of-the-art UAD methods, our method is easier to implement.

Introduction

In the realm of data analysis, anomaly detection (AD) stands as a pivotal challenge with far-reaching implications across various domains, including cybersecurity (Siddiqui et al. 2019; Saeed et al. 2023), healthcare (Yang, Qi, and Zhou 2023; Abououf et al. 2023), finance (Hilal, Gadsden, and Yawney 2022), and industrial processes (Fan, Chow, and Qin 2022; Roth et al. 2022). Existing deep learning-based unsupervised AD methods often rely on an auxiliary learning objective such as auto-encoder, generative model, and contrastive learning. These methods indirectly detect anomalous data using other metrics such as reconstruction error, which lack generalizability and reliability guarantees (Hussain et al. 2023). Explicitly learning a one-class decision boundary may resolve this issue. Many well-known unsupervised AD methods assume the normal training data has a special structure in their data space or embedding space (Schölkopf et al. 2001; Tax and Duin 2004; Ruff et al. 2018; Goyal et al. 2020; Zhang

et al. 2024; Xiao et al. 2025). Such assumptions may not hold or be guaranteed in practice and sometimes place a burden on the model training (Cai and Fan 2022; Fu, Zhang, and Fan 2024). For instance, in deep SVDD (Ruff et al. 2018), the optimal decision boundary in the embedding space may be very different from the learned hypersphere, leading to unsatisfactory detection performance (Zhang et al. 2024).

Given that tabular data is probably the most common data type and other types of data such as images can be converted to tabular data using feature encoders or pretrained models, in this work, we focus on the tabular data only. We propose a novel unsupervised AD method for tabular data without making any assumption about the distribution of normal data. Since hard anomalies are often close to normal data, it is reasonable to hypothesize that hard anomalies are special cases of perturbed samples of normal data. Therefore, if the diversity of the perturbations or added noises on normal data is sufficiently high, we can obtain hard anomalies. Consequently, if a model can recognize highly diverse perturbations or noises, it can detect hard anomalies as well as easy anomalies successfully. By directly learning from the diverse noise patterns and the clean data patterns, we can learn an effective decision boundary around the normal data, generalizing well to unseen data.

Our contributions are highlighted as follows.

- We propose a novel AD method for tabular data using noise evaluation. Our scheme generates highly diverse noise-augmented instances for the normal samples. By evaluating the noise magnitude, our method can accurately identify anomalies.
- The proposed method provides a simple yet effective scheme that does not make assumptions about the normal training data. In addition, the noise generation is straightforward without any extra training. Compared with (Wang et al. 2021; Goyal et al. 2020; Yan et al. 2021; Cai and Fan 2022), our method is more lightweight for training (requires less module for training.)
- We theoretically prove the generalizability and reliability of the proposed method.
- We conduct extensive empirical experiments on 47 real datasets in an unsupervised anomaly detection setting and 25 real-world tabular datasets in a one-class classification setting to demonstrate the performance of the proposed

*Corresponding author.

schemes. The results show that our method achieved superior performance compared with 12 baseline methods including the state-of-the-art.

Related Work

Unsupervised anomaly detection (UAD) is also known as the one-class classification (OCC) problem, in which all or most training samples are assumed to be normal. The learning objective is to learn a decision boundary that distinguishes whether a sample belongs to the same distribution of the normal training data or not. There is another similar problem setting known as outlier detection on contaminated dataset (Huang et al. 2024). The goal is to detect noised samples or outliers within the training data (Ding, Zhao, and Akoglu 2022). This line of work is orthogonal to our UAD problem.

In the past decades, many UAD methods have been studied (Liu, Ting, and Zhou 2008; Chang et al. 2023). Traditional methods like proximity-based (Breunig et al. 2000; Angiulli and Pizzuti 2002; Papadimitriou et al. 2003; He, Xu, and Deng 2003), probability-based (Yang, Latecki, and Pokrajac 2009; Zong et al. 2018; Li et al. 2020), and one-class support vector machine (Schölkopf et al. 2001; Tax and Duin 2004) approaches struggle with high dimensionality and complex data structures. Deep neural network-based methods have been proposed to address these issues. For instance, auto-encoder methods identify outliers by detecting high reconstruction errors, as outlier samples do not conform to historical data patterns (Aggarwal 2016; Chen et al. 2017; Wang et al. 2021). Generative model methods compare latent features or generated samples to spot anomalies (Schlegl et al. 2017; Liu et al. 2019; Zhang et al. 2023; Tur et al. 2023; Xiao et al. 2025). For example, (Xiao et al. 2025) proposed an inverse generative adversarial network that converts the data distribution to a compact Gaussian distribution, based on which the density of test data can be calculated for anomaly detection. Contrastive learning (Sohn et al. 2020; Jin et al. 2021; Shenkar and Wolf 2022) leverages feature representation differences to detect anomalies. Unlike autoencoder-based methods that focus on reducing reconstruction error and reducing dimensionality to remove noise, our method aims to evaluate noise level, which is similar to the denoising diffusion model (Ho, Jain, and Abbeel 2020).

Some works explicitly build an anomaly detection or OCC objective (Ruff et al. 2018; Goyal et al. 2020; Yan et al. 2021; Chen et al. 2022; Cai and Fan 2022). For instance, Deep SVDD (Ruff et al. 2018) trains a neural network to construct a hypersphere in the output space to enclose the normal training data. DROCC (Goyal et al. 2020) assumes normal samples lie on low-dimensional manifolds and treats identifying the ideal hypersphere as an adversarial optimization problem. PLAD (Cai and Fan 2022) outputs an anomaly score by learning a small perturbation of normal data as the negative sample with a classifier. Unlike PLAD, which uses extra additive and multiplicative perturbations requiring a perturbator, our method generates negative samples without extra parameters, making the training efficient.

It is worth mentioning that there are vision UAD methods utilizing synthetic anomalous data. However, the normality and abnormality can be visualized directly and there naturally

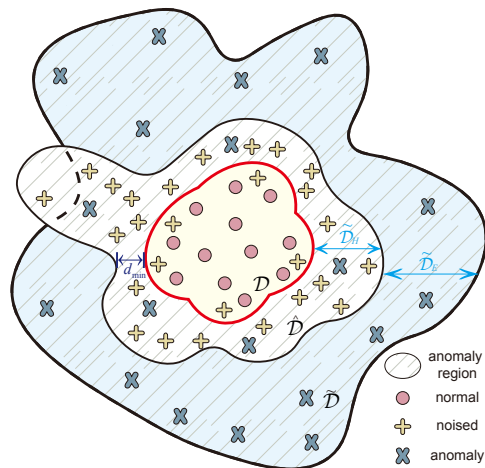


Figure 1: An illustration of the allocation of normal, noised, and true anomalous samples. \mathcal{D} , $\hat{\mathcal{D}}$, and $\tilde{\mathcal{D}}$ are the normal, noised, and anomalous distributions respectively. $\tilde{\mathcal{D}}$ is composed of a hard part $\tilde{\mathcal{D}}_H$ and an easy part $\tilde{\mathcal{D}}_E$. Theorem 1 and Theorem 2 are for $\tilde{\mathcal{D}}_H$ and $\tilde{\mathcal{D}}_E$ respectively.

Symbol	Description	Symbol	Description
x	a scalar	\mathbf{x}_i	a vector with index i
\mathcal{X}_i	a set with index i	$[i]$	the set $\{1, 2, \dots, i\}$
\mathcal{D}	a distribution	$\ \cdot\ _2$	ℓ_2 norm of vector
$\ \cdot\ _1$	ℓ_1 norm of vector	$ \mathbf{x} $	element-wise absolute subset
\cup	union of sets	\subset	subset
H	entropy	h_θ	model parameterized by θ

Table 1: Notations

exists prior knowledge about anomalies (e.g., visible spots or blurs) in visual data. Regarding tabular data, we do not have such prior knowledge. Vision AD methods require prior pre-trained models or external reference datasets to obtain the negative samples, such as DREAM (Zavrtanik, Kristan, and Skočaj 2021) with the Describable Textures Dataset or AnomalyDiffusion (Hu et al. 2024) with a stable diffusion model. Such resources are costly and violate the principle of unsupervised learning. Tabular data, however, spans diverse domains (e.g., medical, industrial), making external datasets and pretrained model infeasible.

Proposed Method

Problem Formulation and Notations

Given a data set $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, where the element \mathbf{x}_i are drawn from an unknown distribution \mathcal{D} in \mathbb{R}^d (deemed as a normal data distribution). In the context of UAD, the primary objective is to develop a function $f : \mathbb{R}^d \rightarrow \{0, 1\}$, which effectively discriminates between in-distribution (normal) and out-of-distribution (anomalous) instances. This discriminative function is formulated to assign a binary label, where $f(x) = 1$ indicates x does not belong to \mathcal{D} and $f(x) = 0$ corresponds to x coming from \mathcal{D} . The main notations used in this paper are shown in Table 1.

Anomalous Data Decomposition

Let $\tilde{\mathcal{X}}$ be the set consisting of all anomalous data drawn from some unknown distribution $\tilde{\mathcal{D}}$ deemed as an anomalous distribution. For any $\tilde{x} \in \tilde{\mathcal{X}}$, we decompose it as

$$\tilde{x} = x + \epsilon, \quad (1)$$

where $x \in \mathcal{D}$ is a normal counterpart of \tilde{x} and ϵ denotes the derivation of \tilde{x} from x . The magnitude of ϵ , denoted as $\|\epsilon\|_1$, measures how anomalous \tilde{x} is. Note that this decomposition is not unique and hence one may seek the one with the smallest $\|\epsilon\|_1$. If we can learn a model h to predict ϵ for \tilde{x} , i.e.,

$$\epsilon = h(\tilde{x}), \quad (2)$$

we will be able to determine whether \tilde{x} is normal or not according to $\|\epsilon\|_1$. The challenge is that there is no available information about $\tilde{\mathcal{X}}$ in the training stage and we can only utilize \mathcal{X} .

Although $\tilde{\mathcal{X}}$ is unknown, we further theoretically partition $\tilde{\mathcal{X}}$ into two subsets without overlap, i.e.,

$$\tilde{\mathcal{X}} \triangleq \tilde{\mathcal{X}}_E \cup \tilde{\mathcal{X}}_H, \quad \tilde{\mathcal{X}}_E \cap \tilde{\mathcal{X}}_H = \emptyset, \quad \tilde{\mathcal{X}}_E \sim \tilde{\mathcal{D}}_E, \quad \tilde{\mathcal{X}}_H \sim \tilde{\mathcal{D}}_H. \quad (3)$$

$\tilde{\mathcal{X}}_E$ denotes an easy set, drawn from the easy part $\tilde{\mathcal{D}}_E$ of $\tilde{\mathcal{D}}$, in which $\|\epsilon\|_1$ for each sample is sufficiently large, while $\tilde{\mathcal{X}}_H$ denotes a hard set, drawn from the hard part $\tilde{\mathcal{D}}_H$ of $\tilde{\mathcal{D}}$, in which $\|\epsilon\|_1$ for each sample is small. After the partition, we can assert that $\tilde{\mathcal{X}}_H$ is closer to the normal data. Consequently, it is easier for a model to recognize samples in $\tilde{\mathcal{X}}_E$ than those in $\tilde{\mathcal{X}}_H$, as the $\|\epsilon\|_1$ values of samples in $\tilde{\mathcal{D}}_E$ are significantly larger than those in $\tilde{\mathcal{D}}_H$. Figure 1 provides an intuitive example. The ultimate goal is to learn the decision boundary around the normal data.

Here, we focus on how to detect the samples in $\tilde{\mathcal{X}}_H$ or drawn from $\tilde{\mathcal{D}}_H$. Since $\tilde{\mathcal{X}}_H$ is very close to the normal data, it is reasonable to hypothesize that hard anomalies are special cases of perturbed samples of normal data. We propose to generate a noisy dataset $\hat{\mathcal{X}} \subset \mathbb{R}^d$ from \mathcal{X} by adding various noise to \mathcal{X} , and assume $\hat{\mathcal{X}}$ is drawn from certain perturbed distribution $\hat{\mathcal{D}}$, i.e.,

$$\hat{\mathcal{X}} \leftarrow \text{Gen}(\mathcal{X}) \sim \hat{\mathcal{D}}, \quad (4)$$

where Gen denotes the noisy data generator and $|\hat{\mathcal{X}}| \gg N$. Let the diversity of added noise is sufficiently large, such that

$$\tilde{\mathcal{X}}_H \subset \hat{\mathcal{X}}, \quad (5)$$

i.e., anomaly patterns of the hard set $\tilde{\mathcal{X}}_H$ are included in $\hat{\mathcal{X}}$. Even if (5) does not hold, as shown by Theorem 1, it is still possible to obtain correct detection, provided that $\hat{\mathcal{D}}$ is not too far from $\tilde{\mathcal{D}}_H$, i.e.,

$$\text{dist}(\hat{\mathcal{D}}, \tilde{\mathcal{D}}_H) \leq \gamma, \quad (6)$$

where $\text{dist}(\cdot, \cdot)$ is some distance or divergence measure between two distributions and γ is not too large. Therefore, a model h learned from $\hat{\mathcal{X}}$ is able to generalize to $\tilde{\mathcal{X}}_H$ and then detect anomaly.

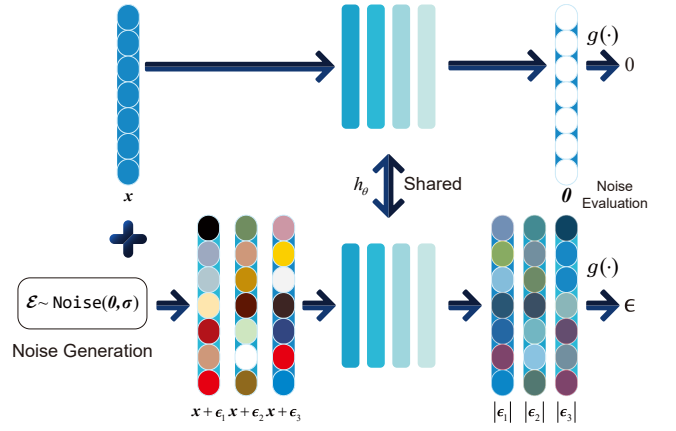


Figure 2: The training process of noise evaluation model. Noise with 0 mean and σ standard deviation is added to the original data x to create noised versions $\hat{x} = x + \epsilon$. The model h_θ is trained to discern the zero vector for the original data and identify the noise vector $|\epsilon|$ for the noised data. The final anomaly decision is made using an aggregation function $g(\cdot)$, where high-magnitude noise indicates abnormality.

Noise Evaluation Model

To generate $\hat{\mathcal{X}}$, we add random noise to the elements of each sample $x \in \mathcal{X}$. This operation will reduce the quality of the data, that is, the quality of $\hat{\mathcal{X}}$ is lower than that of \mathcal{X} , supported by

Proposition 1. *Adding random noises independently to the entries of \mathcal{X} makes the data more disordered (higher entropy).*

We would like to learn a deep neural network $h(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ parameterized by θ to quantify the quality of the input data. Its output is a vector with the same size as the input while each entry tells whether the input feature is noised or not. We generate the noised dataset by

$$\hat{\mathcal{X}} = \hat{\mathcal{X}}_1 \cup \hat{\mathcal{X}}_2 \cdots \cup \hat{\mathcal{X}}_K, \quad (7)$$

where each $\hat{\mathcal{X}}_k$ is composed of the samples generated by

$$\hat{x} = x + \epsilon, \quad \epsilon \sim \text{Noise}_k(\mathbf{0}, \sigma_k), \quad x \in \mathcal{X}. \quad (8)$$

In (8), $\text{Noise}_k(\mathbf{0}, \sigma_k)$ (to be detailed in the next section) is a multivariate noise distribution with $\mathbf{0}$ mean value and σ_k standard deviation in \mathbb{R}^d , abbreviated as $\text{Noise}(\sigma_k)$. For instance, Noise_k can be a Gaussian distribution. We also denote the standard deviation of the noise as *noise level*¹. According to (7) and (8), we see that $\hat{\mathcal{X}}$ can contain different types of noise distributions with different standard deviations, and the diversity is controlled by $\sum_{k=1}^K \|\sigma_k\|$. Refer to Table 3 for the choices of the noise generator. Then, we write the learning objective as

$$\min_{\theta} \sum_{x_i \in \mathcal{X}} \|h_\theta(x_i) - \mathbf{0}\|_2^2 + \sum_{\hat{x}_i \in \hat{\mathcal{X}}} \|h_\theta(\hat{x}_i) - |\epsilon_i|\|_2^2, \quad (9)$$

where $\mathbf{0}$ is a d -dimension vector with all zero values, ϵ is drawn from some $\text{Noise}(0, \sigma)$, and $\hat{x} = x + \epsilon$ for some

¹We use noise level and standard deviation of noise interchangeably in this paper.

$\mathbf{x} \in \mathcal{X}$. Note that in (9), there is no hyperparameter to determine except the network structure, while many previous methods such as (Ruff et al. 2018; Goyal et al. 2020; Cai and Fan 2022; Zhang et al. 2024) have at least one more crucial hyperparameter to tune in the learning objective.

In (9), the absolute value is used to let the model output a positive signal indicating an anomaly. Here, we only require the model to predict how much the noise is instead of the exact noise value, which also reduces the difficulty of the learning process. In addition, our method is closely related to the denoising score matching (Song and Ermon 2019; Vincent 2011), which estimates the gradient of the probability density of the data distribution. We prove (9) is a lower bound of denoising score matching learning objective in Appendix A. Hence, our method also learns data distribution \mathcal{D} by proxy. Predicting the noise magnitude has the following advantages.

Claim 1. *Estimating the element-wise magnitude of noise enables the model to quantify the deviation of a sample from the normal data distribution \mathcal{D} . It improves the ability to identify specific features with abnormalities.*

Since the model output is a d -dimension vector, we introduce an aggregation operation after the training, $g(\cdot)$, to map the output vector to a scalar for anomaly detection. The operation in our design can be maximum, minimum, mean, median, or a combination. The final anomaly score is determined by

$$\text{score}(\mathbf{x}) := g(h_{\theta}(\mathbf{x})). \quad (10)$$

Taking the maximum as an example, we have $\text{score}(\mathbf{x}) = \max_{i \in [d]} [h_{\theta}(\mathbf{x})]_i$. We need to determine a threshold $\tau > 0$ for the anomaly score. If $\text{score}(\mathbf{x}) > \tau$, \mathbf{x} is anomalous. An overview of our model is shown in Figure 2.

Noise Generation and Model Training

The noise generation in our design is not arbitrary. To cover the sampling space of the noised distributions as much as possible, we randomly generate the noise for each training sample in terms of noise level and position. An intuitive example is as follows:

$$\boldsymbol{\epsilon} = \left[\begin{array}{ccccccc} \underbrace{\epsilon_1, \epsilon_2}_{\sim \text{Noise}(\sigma_1)} & \underbrace{\epsilon_3, \epsilon_4, \dots, \epsilon_i, \epsilon_{i+1}}_{\sim \text{Noise}(\sigma_2)} & \underbrace{\epsilon_{i+2}, \dots, \epsilon_{d-2}}_{\sim \text{Noise}(\sigma_3)} & \underbrace{\epsilon_{d-1}, \epsilon_d}_{\sim \text{Noise}(\sigma_1)} \\ \underbrace{\epsilon_1, \epsilon_2}_{\sim \text{Noise}(\sigma_1)} & \underbrace{\epsilon_3, \epsilon_4, \dots, \epsilon_i, \epsilon_{i+1}}_{\sim \text{Noise}(\sigma_2)} & \underbrace{\epsilon_{i+2}, \dots, \epsilon_{d-2}}_{\sim \text{Noise}(\sigma_3)} & \underbrace{\epsilon_{d-1}, \epsilon_d}_{\sim \text{Noise}(\sigma_1)} \end{array} \right]. \quad (11)$$

For one type of distribution (e.g. Gaussian), there are often two hyper-parameters to control the noise generation process. The first one is σ_{max} , denoting the maximum noise level. The other is m , describing how many parts of the feature vector are added by the same level of noise. Therefore, for an \mathbf{x} , different elements may be corrupted by different levels of noise, and, if necessary, one \mathbf{x} can produce multiple $\hat{\mathbf{x}}$ with different types of distribution. A detailed process for generating noised samples on a batch of data is in Algorithm 1. According to Algorithm 1, the noise generation time complexity is $\mathcal{O}(bd)$. In contrast, other methods involving perturbation (Cai and Fan 2022; Qiu et al. 2021) and adversarial sample (Goyal et al. 2020) have time complexity with $\mathcal{O}(bdW)$, where W is workload related to a neural network module. Compared with them, our noise generation is more efficient, where no

Algorithm 1: Noise Generation

Input: maximum noise level σ_{max} , number of noise distributions m , batch size b , the dimensionality of data d .

- 1: $\Delta \leftarrow [0, \frac{1}{m}\sigma_{max}, \frac{2}{m}\sigma_{max}, \dots, \frac{m}{m}\sigma_{max}]$ \triangleright make m intervals
- 2: Initialize \mathcal{E} with empty
- 3: **for** $j \in \{1, \dots, b\}$ **do**
- 4: **for** $i \in \{1, \dots, m\}$ **do**
- 5: //random noise level
- 6: $\hat{\sigma} \leftarrow \text{Uniform}(\Delta[i], \Delta[i+1])$
- 7: $\epsilon_{\lfloor \frac{i-1}{m}d : \frac{i}{m}d \rfloor} \leftarrow \text{Noise}_{\mathcal{E}}(0, \hat{\sigma})$ //generate noise from m noise levels for m parts
- 8: **end for**
- 9: //shuffle position of noise elements
- 10: $\boldsymbol{\epsilon} \leftarrow \text{Shuffle}(\boldsymbol{\epsilon})$
- 11: $\mathcal{E}[j] \leftarrow \boldsymbol{\epsilon}$
- 12: **end for**

Output: Generated noise \mathcal{E} for a batch of input data

learnable parameter is required. A comparative study on time cost is shown in Appendix I.

In Figure 2, the lower pathway illustrates the noise synthesis mechanism, where a noise vector $\boldsymbol{\epsilon}$ (mean 0, standard deviation $\boldsymbol{\sigma}$) is randomly generated and added to the input \mathbf{x} , producing a noise-augmented variant $\hat{\mathbf{x}} = \mathbf{x} + \boldsymbol{\epsilon}$. Multiple noised samples can be generated from a single input. Both \mathbf{x} and $\hat{\mathbf{x}}$ are processed by the noise evaluation network h_{θ} , optimized to regress towards zero for \mathbf{x} and estimate the noise vector $|\boldsymbol{\epsilon}|$ for $\hat{\mathbf{x}}$. Training details are in Appendix B, Algorithm 2. Notice that for each training epoch, we randomly generate a new noised instance for each training sample. This helps us enlarge the sampling number from the noise distribution, and avoid over-fitting on some ineffective noise. There are some optional noise generation schemes such as using different noise types, different noise levels, and different noise ratios, which are studied later. These options add several hyper-parameters to our methods. However, it is still simpler than many recent methods. We compare them in Appendix J.

Theoretical Analysis

In the proposed method, we learn a one-class classification decision boundary closely around the normal samples. In Figure 1, we divide the anomaly region into two non-overlap distributions, easy $\hat{\mathcal{D}}_E$, and hard $\hat{\mathcal{D}}_H$, respectively. In this section, we theoretically show the anomaly detection ability of the model meeting easy anomaly samples \mathcal{X}_E from $\hat{\mathcal{D}}_E$ and hard anomaly samples \mathcal{X}_H from $\hat{\mathcal{D}}_H$ in Theorem 1 and Theorem 2, respectively.

Without loss of generality, we let $h_{\theta} \in \mathcal{H}$, where \mathcal{H} is the hypothesis space of ReLU-activated neural network with L layers and the number of neurons at each layer is in the order of p . We give the following definition.

Definition 1. *For the noise evaluation hypothesis h_{θ} , the risk of the hypothesis on a distribution \mathcal{D} is defined by*

the probability according to \mathcal{D} that a processed hypothesis $I(g(h_\theta(\mathbf{x}))) > \tau$ disagrees with a labeling function $\rho : \mathbb{R}^d \rightarrow \{0, 1\}$. Mathematically,

$$\varepsilon_{\mathcal{D}}(h) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [|I(g(h_\theta(\mathbf{x}))) > \tau| - \rho(\mathbf{x}) |],$$

where $I(\cdot)$ is the indicator function, $\tau > 0$ is some threshold, and $\rho(\cdot)$ is the true labeling function.

Definition 1 is a form of the disagreement metric (Hanneke et al. 2014). Based on Definition 1 and results of (Ben-David et al. 2010; Bartlett et al. 2019), we can provide the following theoretical guarantee (proved in Appendix C) for the hard anomalous data $\tilde{\mathcal{X}}_H$.

Theorem 1. (Generalization Error Bound) Let $\varepsilon_{\tilde{\mathcal{D}}_H}(h)$ and $\varepsilon_{\hat{\mathcal{D}}}(h)$ be the risks of hypothesis h on $\tilde{\mathcal{D}}_H$ and $\hat{\mathcal{D}}$, respectively. If $\hat{\mathcal{X}}$ and $\tilde{\mathcal{X}}_H$ are unlabeled samples of size N each, drawn from $\hat{\mathcal{D}}$ and $\tilde{\mathcal{D}}_H$ respectively, then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, for every $h \in \mathcal{H}$:

$$\varepsilon_{\tilde{\mathcal{D}}_H}(h) \leq \varepsilon_{\hat{\mathcal{D}}}(h) + \frac{1}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{X}}_H, \hat{\mathcal{X}}) + 4\sqrt{\frac{2d_{vc} \log(2N) + \log(\frac{2}{\delta})}{N}} + \lambda,$$

where $d_{vc} = \mathcal{O}(pL \log(pL))$. $\lambda = \varepsilon_{\hat{\mathcal{D}}}(h^*) + \varepsilon_{\tilde{\mathcal{D}}_H}(h^*)$ where h^* is the ideal joint hypothesis that minimize $\varepsilon_{\hat{\mathcal{D}}}(h) + \varepsilon_{\tilde{\mathcal{D}}_H}(h)$.

The definition of the divergence $\hat{d}_{\mathcal{H}\Delta\mathcal{H}}$ is shown in Appendix C for simplicity. As shown in Figure 1, Theorem 1 describes that if the divergence $\hat{d}_{\mathcal{H}\Delta\mathcal{H}}$ between the perturbed training data $\hat{\mathcal{X}}$ and the hard test data $\tilde{\mathcal{X}}_H$ is small, the model can correctly identify the anomaly. In other words, the generated data $\hat{\mathcal{X}}$ can be very useful for learning a reasonable detection model h_θ . Since there is no available information about $\tilde{\mathcal{X}}_H$ during the training, we cannot obtain the divergence $\hat{d}_{\mathcal{H}\Delta\mathcal{H}}$ in real practice. Hence, we always standardize the normal training data and make the added noise level relatively small. In the ablation study, we show that a large noise level is harmful to the training. Note that this guarantee does not apply to $\tilde{\mathcal{X}}_E$ because $\tilde{\mathcal{D}}_E$ may be very far from $\hat{\mathcal{D}}$. The following context will provide a guarantee for detecting $\tilde{\mathcal{X}}_E$. We make the following assumption and present the theoretical result (proved in Appendix C).

Assumption 1. For any \mathbf{x} and $\tilde{\mathbf{x}}$ drawn from \mathcal{D} and $\tilde{\mathcal{D}}_E$ respectively, there exists a constant $c > 0$ such that $c \|\mathbf{x} - \tilde{\mathbf{x}}\| \leq |g(h_\theta(\mathbf{x})) - g(h_\theta(\tilde{\mathbf{x}}))|$.

Theorem 2. Let $d_{\min} = \inf_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}_E} \|\mathbf{x} - \tilde{\mathbf{x}}\|$ (shown by Figure 1). Suppose $g(h_\theta(\mathbf{x})) < \epsilon$ for any $\mathbf{x} \in \mathcal{D}$, where $\epsilon \geq 0$. Then, under Assumption 1, any anomalous samples drawn from $\tilde{\mathcal{D}}_E$ can be successfully detected if $d_{\min} > \max\{\epsilon/c, \tau/c\}$.

This theorem provides a theoretical guarantee for our method to detect any anomalous sample in $\tilde{\mathcal{X}}_E$ or drawn from $\tilde{\mathcal{D}}_E$ (depicted by the blue double-headed arrow in Figure 1) as anomaly successfully. When c is larger, the detection is easier. We can derive a bound for c with more assumptions: if the learned model h is bijective and L -Lipschitz, then h^{-1} is $1/L$ -Lipschitz, meaning that $c = \alpha L$, where

$\alpha = \inf_{z \in \text{dom}(g)} \|\nabla g(z)\|$. Using an invertible neural network (Behrmann et al. 2019) could achieve this. For an h with q layers, spectral norm $\|\mathbf{W}_i\|_\sigma$ for layer i , and ρ -Lipschitz activation, we have $c = \alpha \prod_{i=1}^q \rho^q \|\mathbf{W}_i\|_\sigma^q$. In our experiments, we used both MLP and ResMLP without dimensionality reduction. According to Theorem 1 of (Behrmann et al. 2019), ResMLP is invertible if each layer’s Lipschitz constant is under 1, which is easy to ensure.

To sum up, Theorem 1 and Theorem 2 provide guarantees for detecting \mathcal{X}_H and \mathcal{X}_E respectively. Therefore, our method is theoretically guaranteed to detect $\tilde{\mathcal{X}}$, although we never use any real anomalous data in the training stage.

Experiments

Experimental Settings

Datasets We evaluate our method in two common settings: unsupervised anomaly detection and one-class classification. In the anomaly detection setting, where anomalous samples are few, we use 47 real-world tabular datasets² from (Han et al. 2022), covering domains like healthcare, image processing, and finance. For one-class classification setting, we collected 25 benchmark tabular datasets used in previous works (Pang et al. 2021; Shenkar and Wolf 2022). The raw data was sourced from the UCI Machine Learning Repository (Kelly, Longjohn, and Nottingham) and their official websites. For categoric value, we use a one-hot encoding. We test on all classes in multi-class datasets, reporting the average performance score per class, which is similar to one-class classification on image dataset (Cai and Fan 2022). For datasets with validation/testing sets, we train on all normal samples. If only a training set is available, we randomly split 50% of normal samples for training and use the rest with anomalous data for testing. Data is standardized using the training set’s mean and standard deviation. Dataset details are in Table 4 of Appendix D.

Baseline Methods We select 12 baseline methods for comparative analysis, including probabilistic-based, proximity-based, deep neural network-based, ensemble-based methods, and recent UAD methods that can be applied to tabular data. They are Isolated Forest (**IForest**) (Liu, Ting, and Zhou 2008), **COPOD** (Li et al. 2020), auto-encoder (**AE**) (Aggarwal 2016), **KNN** (Angiulli and Pizzuti 2002), Local Outlier Factor (**LOF**) (Breunig et al. 2000), **DeepSVDD** (Ruff et al. 2018), **AnoGAN** (Schlegl et al. 2017), **ECOD** (Li et al. 2023), **SCAD** (Shenkar and Wolf 2022), **NeuTraLAD** (Qiu et al. 2021), **PLAD** (Cai and Fan 2022), and **DPAD** (Fu, Zhang, and Fan 2024). For DPAD, PLAD, SCAD, and NeuTraLAD, we use the code provided by the authors. For the other baseline methods, we utilize PyOD, a Python library developed by (Zhao, Nasrullah, and Li 2019). The default settings are adopted. We repeat 10 times for each baseline.

Implementation All experiments are implemented by Pytorch (Paszke et al. 2017) on NVIDIA Tesla V100 and Intel Xeon Gold 6200 platform. We utilize two network architectures for the evaluation, VanillaMLP (**MLP**) and **ResMLP**.

²<https://github.com/Minqi824/ADBench/>

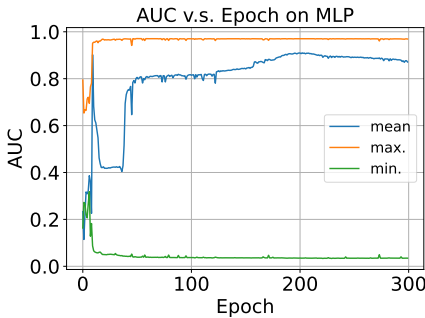


Figure 3: Comparison of different $g(\cdot)$, i.e. mean, maximum, and minimum, on KDD-CUP99, at each optimization epoch.

VanillaMLP is a plain ReLU-activated feed-forward neural network with 4 fully connected layers. ResMLP has a similar architecture to (Touvron et al. 2022) but has a simpler structure. Detailed architectures are in Appendix E. The network model is optimized by AMSGrad (Reddi, Kale, and Kumar 2018) with 10^{-4} learning rate and 5×10^{-4} weight decay. In the results, we adopt Gaussian noise in the noise generation, maximum noise level $\sigma_{max} = 2$, and the number of different noise distributions $m = 3$. To enlarge the noised sample number, we generate 3 noise-augmented instances for the same input instance with 3 different noise ratios, 0.5, 0.8, and 1.0 at each epoch. Unless specified, we train the model for 500 epochs and manually decay the learning rate at the 100-th epoch in a factor of 0.1.

Performance Metric Since the output of h is a noise evaluation vector, an operation $g(\cdot)$ is introduced, as explained by (10). Here we select the maximum. A comparative result on KDD-CUP99 among maximum, minimum, and mean is shown in Figure 3 where we train the model for 300 epochs. The maximum reaches the highest performance and the fastest convergence speed. For the performance metric, we report the area under the ROC curve (AUC) and F1 score. The calculation of the F1 score and the determination of anomaly threshold τ are consistent with (Shenkar and Wolf 2022; Qiu et al. 2021).

Unsupervised Anomaly Detection Results

The average result under the unsupervised anomaly detection setting of ten runs is reported in Figure 4. We conducted a comparative analysis of 11 tabular UAD methods across 47 datasets, evaluating average AUC, average F1, average ranking in AUC, and average ranking in F1. It is seen that our methods not only significantly outperform the AE methods, but also reach the highest ranking. The detailed results on each dataset and p-value from paired t-test are reported in Appendix F (Tables 6 and 7), which emphasizes the statistical significance of the improvements achieved by our methods.

One-Class Classification Results

We compared our noise evaluation method with 12 baseline methods on the OCC dataset setting. The results in Table 2 show that our anomaly detection techniques, Ours-ResMLP and Ours-MLP, consistently outperform baselines across vari-

ous tabular datasets, with mean AUC values of 92.68 ± 11.10 and 92.27 ± 11.1 , indicating greater effectiveness and lower variability. While traditional methods like IForest and KNN perform well, our methods excel, especially on complex datasets like musk and optdigits. Though our methods may not always lead in AUC, the difference is minimal, around 1%. For faster inference, the MLP model is recommended. The last two lines in the table are the p-values of paired t-test of our two methods (ResMLP and MLP) against each baseline method. The paired t-test is based on 25 pairs (as there are 25 datasets). Our approach also achieved the highest average F1 score and rank, 94.18% and 1.52, as detailed in Appendix G.

Compared with many popular UAD methods, our noise evaluation method has the following advantages:

- Noise evaluation shows superior performance in the extensive empirical experiments, indicating our method can accurately identify anomalies in tabular data.
- The generalization and anomaly detection ability of noise evaluation can be theoretically guaranteed, whereas many deep learning based methods lack such guarantees (Husain et al. 2023).
- The implementation of our method is easier. The perturbed sample generation can be pre-generated without extra training. In addition, no hyper-parameter is tuned in the learning objective.

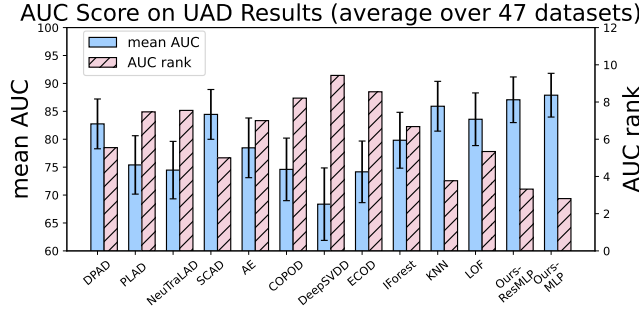
Type	Parameter	Value	Offset
Gaussian	μ, σ	$\mu = 0, \sigma = \hat{\sigma}$	0
Laplace	μ, σ	$\mu = 0, \sigma = \hat{\sigma}$	0
Uniform	a, b	$a = -\sqrt{3}\hat{\sigma}, b = \sqrt{3}\hat{\sigma}$	0
Rayleigh	s	$s = \sqrt{\frac{2}{4-\pi}}\hat{\sigma}$	$-\sqrt{\frac{\pi}{4-\pi}}\hat{\sigma}$
Gamma	α, β	$\alpha = \sqrt{\frac{1}{\beta}}\hat{\sigma}, \beta = \beta$	$-\sqrt{\beta}\hat{\sigma}$
Poisson	λ	$\lambda = \hat{\sigma}$	$-\hat{\sigma}$

Table 3: Eight different noise types are adopted. We adjust the parameter to make it 0 mean value and $\hat{\sigma}$ standard deviation, where $\hat{\sigma}$ is the designated noise level. If the noise is non-negative by default, we offset its mean to 0. For the Gamma noise, β is a non-negative hyper-parameter, where we evaluate $\beta = 1$ and $\beta = 3$. For Salt&Pepper and Bernoulli noise, we generate a probability vector based on a uniform distribution. Then, generate a binary vector using the probability vector. If there is noise, we change the value into the maximum or minimum value in the batch or flip the sign of the element of \mathbf{x} .

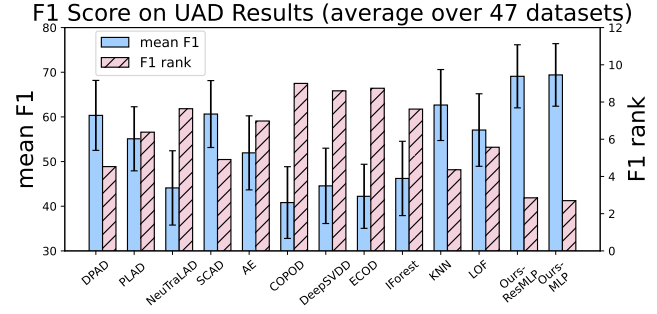
Ablation Study

As discussed in the Proposed Method Section, we have several alternatives. In this section, we study how different noise levels, noise ratios, and noise types affect the performance of our method. We utilize the OCC setting to perform the ablation study. Detailed results are shown in Appendix H.

Sensitivity of Different Noise Levels To explore how different noise levels affect performance, we use Gaussian



(a) AUC and Rank on UAD Results.



(b) F1 and Rank on UAD Results.

Figure 4: AUC (%) and F1 (%) score of the proposed method compared with 11 baselines on 47 benchmark datasets. Each experiment is repeated 10 times with random seed from 0 to 9, and mean value and 95% confidence interval are reported. Rank (the lower the better) is calculated out of 12 tested methods.

Methods	DPAD (2024)	PLAD (2022)	NeuTraLAD (2021)	SCAD (2022)	AE	AnoGAN	COPOD	DeepSVDD	ECOD (2023)	iForest	KNN	LOF	Ours-ResMLP	Ours-MLP
abalone	70.64±13.2	64.23±15.4	45.29±8.8	63.90±14.7	71.69 ±11.1	68.65±14.7	43.57±15.0	59.78±16.1	45.00±16.9	71.55±12.2	69.11±14.1	68.89±13.0	76.90 ±11.0	76.78 ±10.9
arrhythmia	76.38±1.0	63.15±5.1	52.13±0.0	69.60±2.0	76.12±1.8	69.61±2.6	75.23±1.4	72.82±2.2	75.21±1.5	77.06 ±1.8	75.87±2.0	75.85±1.9	77.34 ±2.4	76.77±2.1
breastw	98.36±0.3	87.4±15.0	78.77±3.0	97.29±0.7	98.83±0.4	99.11±0.3	99.29 ±0.2	95.58±1.1	98.62±0.3	99.36 ±0.2	98.84±0.3	95.43±1.6	98.49±0.7	98.84±0.4
cardio	82.13±3.3	74.82±8.9	55.84±2.6	69.37±1.8	83.36 ±0.8	76.78±12.2	65.98±0.7	53.12±7.7	78.03±0.6	81.74±2.0	73.96±1.0	77.37±1.5	85.59 ±2.4	86.01 ±3.3
ecoli	92.25±4.4	68.40±2.6	48.46±12.0	88.17±5.7	93.18±4.0	90.86±6.2	49.16±22.2	87.63±5.8	50.98±11.6	93.21 ±3.6	93.19±4.2	91.42±5.7	92.98±3.7	93.53 ±3.7
glass	76.73±7.4	64.8±10.5	57.68±9.0	78.89±6.0	73.39±10.4	74.76±11.0	44.61±26.1	75.39±8.2	44.75±24.2	76.96±10.5	76.66±9.0	71.09±9.7	83.12 ±10.1	84.41 ±10.2
ionosphere	96.66±0.4	64.26±15.2	90.77±2.4	96.83±0.7	90.40±1.2	86.19±5.9	80.05±1.6	95.56±1.1	74.33±1.6	89.53±2.9	97.47 ±0.8	95.44±1.3	97.53 ±0.3	97.11±0.7
kdd	81.76±11.9	90.87±6.6	93.6±1.7	91.56±4.0	95.12±3.3	92.71±2.5	76.51±1.2	72.81±27.8	78.75±1.3	96.12 ±0.3	94.45±0.2	85.05±0.2	95.69 ±1.2	96.92 ±0.6
letter	93.39±3.3	59.05±15.3	90.95±4.6	99.26 ±0.5	88.91±4.3	80.23±8.5	50.03±15.8	94.25±2.4	50.09±15.4	91.15±3.4	97.95±1.6	95.61±2.6	99.42 ±0.3	99.42 ±0.4
lympho	74.21±3.6	64.96±8.1	46.6±2.2	74.91±4.6	75.84±3.7	70.15±6.7	53.38±3.7	73.37±5.8	53.52±3.5	75.14±5.3	77.86 ±3.2	76.50±3.2	81.14 ±4.7	82.37 ±5.0
mammo.	88.32±1.7	82.38±2.3	65.57±2.3	78.18±2.2	85.65±1.1	78.67±13.4	90.54 ±0.1	69.87±7.1	90.63 ±0.1	87.91±0.7	87.55±0.3	84.12±1.2	90.11±0.9	89.89±0.4
mulcross	100.0 ±0.0	99.90±0.0	76.59±11.3	100.0 ±0.0	100.0 ±0.0	99.89±0.1	93.24±0.0	100.0 ±0.0	95.97±0.1	99.90±0.1	100.0 ±0.0	100.0 ±0.0	100.0 ±0.0	100.0 ±0.0
musk	71.36±2.9	74.42±3.3	78.86±1.6	81.30±0.7	30.68±0.6	41.37±12.2	32.03±0.4	62.02±6.0	30.16±0.4	46.74±3.7	81.94 ±0.4	81.31±0.6	85.06 ±0.5	83.69 ±1.0
optdigits	93.43±4.3	80.43±15.5	67.3±15.0	97.80 ±1.8	95.69±2.6	90.04±6.6	71.72±8.8	94.86±3.7	68.69±10.1	97.55±1.7	97.25±2.2	97.17±2.3	97.93 ±1.7	98.00 ±1.7
pendigits	98.04±2.1	86.45±14.9	97.23±2.4	99.65 ±0.3	95.46±4.7	92.56±9.8	49.90±25.2	94.25±5.3	49.93±22.9	98.08±1.7	99.59±0.5	99.11±1.0	99.82 ±0.2	99.84 ±0.2
pima	70.5±1.7	63.67±5.7	48.63±4.6	67.97±1.3	69.73±1.5	70.94±5.0	65.58±1.0	58.18±2.9	59.46±1.3	73.33±1.4	74.63 ±1.0	71.06±1.5	72.42 ±2.6	73.53 ±2.9
satimage	92.83±5.1	54.83±19.4	47.98±12.0	93.35±4.8	93.39±5.4	91.98±6.8	69.24±26.8	75.10±7.7	62.75±22.8	95.33 ±3.8	94.85 ±4.4	89.67±6.6	94.51±3.6	93.72±4.7
seismic	74.81 ±0.8	72.59±2.8	67.75±1.2	72.47±0.8	71.13±0.5	72.82±1.5	73.87±0.5	54.32±10.2	70.17±0.4	73.48±0.6	74.57 ±0.5	60.87±1.6	71.38±1.3	72.62±1.9
shuttle	98.93 ±1.4	88.55±14.2	97.42±2.4	98.30±2.0	92.49±8.8	88.42±16.2	31.08±31.4	96.51±6.5	35.53±30.8	91.55±6.7	98.66±1.4	95.60±6.2	99.35 ±0.8	99.40 ±1.0
speech	54.38±2.7	57.32±4.9	56.31±4.8	57.85 ±2.4	47.08±0.5	50.33±5.4	49.15±0.6	52.20±3.8	47.08±0.5	46.88±1.6	48.67±0.7	49.81±0.5	57.14±2.6	58.38 ±2.1
thyroid	91.96±1.0	97.02 ±1.7	80.02±3.4	86.53±1.4	83.38±0.6	72.63±4.3	78.46±0.0	57.35±3.5	79.17±0.0	90.08±1.4	91.71±0.0	88.02±0.0	97.47 ±0.1	96.91 ±0.2
vertebral	87.31±2.6	56.47±28.6	50.35±3.5	78.56±3.3	88.58±1.3	87.96±3.2	78.63±2.0	80.43±3.8	60.02±2.9	86.08±2.3	87.18±1.0	88.20±1.5	90.91 ±1.4	89.91 ±1.2
vowels	80.22±4.1	79.77±17.1	97.25±1.1	99.52 ±0.2	62.80±0.8	59.32±7.8	49.73±0.8	72.01±6.5	59.39±0.7	78.17±2.2	97.16±0.4	95.53±0.7	99.42 ±0.3	99.14 ±0.2
wbc	95.10±1.1	68.07±15.0	53.01±12.6	93.11±1.3	95.63±0.6	94.01±2.5	86.53±1.0	90.22±3.3	62.46±1.4	95.74±0.8	94.68±0.9	95.09±0.9	96.98 ±1.1	96.73±1.1
wine	95.68±4.5	67.13±7.8	43.95±12.0	86.21±7.8	97.00±4.1	94.47±10.4	49.37±7.8	95.46±5.5	49.37±11.5	96.23±3.1	97.81±2.8	97.98±2.4	98.76 ±1.8	98.33 ±2.3
mean	88.05±12.4	67.42±19.6	71.45±22.6	88.59±15.3	85.68±13.2	80.17±14.8	53.38±22.4	80.85±19.1	52.77±20.9	87.0±12.8	90.37 ±13.5	88.75±13.4	92.68 ±11.0	92.27 ±11.1
mean rank	4.96	9.16	10.16	5.96	6.00	8.12	10.12	8.68	10.40	4.96	3.92	5.88	2.04	1.68
p-value (ResMLP)	0.0003	0.0000	0.00000	0.0001	0.0066	0.0002	0.0000	0.0000	0.0000	0.0054	0.0025	0.0000	-	0.47
p-value (MLP)	0.0002	0.0000	0.0000	0.0001	0.0051	0.0002	0.0000	0.0000	0.0000	0.0036	0.0022	0.0000	0.47	-

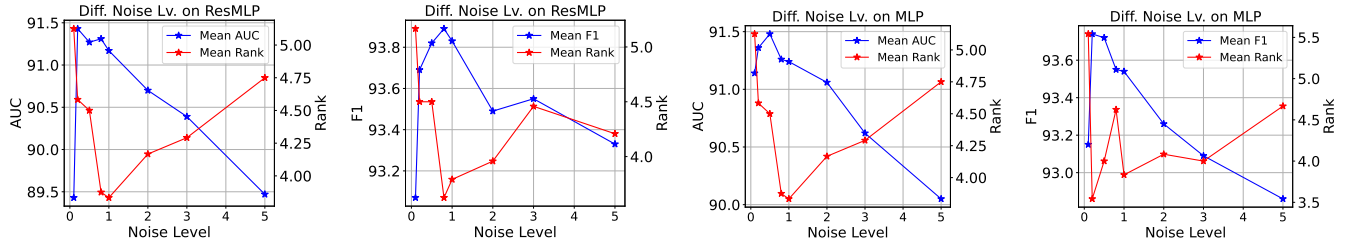
Table 2: AUC (%) score of the proposed method compared with 12 baselines on 25 benchmark datasets. Each experiment is repeated 10 times with random seed from 0 to 9, and mean value and standard deviation are reported. Mean is the average AUC score under all experiments. Mean rank (the lower the better) is calculated out of 10 tested methods. Mammo. refers to mammography.

noise and generate 3 noised instances per training batch, consistent with previous settings. We test noise levels in [0.1, 0.2, 0.5, 0.8, 1.0, 2.0, 3.0, 5.0]. The results are reported in Figure 5. Results show that too small noise levels confuse the model due to the minimal distance between normal samples and anomalies, while too high noise levels expand the output value range and sampling space, reducing effectiveness as theorem 1 suggested. Hence, the optimal noise level is around 1.0.

Sensitivity of Different Noise Types We explore different noise types with a mean of 0 and a standard deviation (noise level) of σ . We use Salt&Pepper noise, Gaussian, Laplace, Uniform, Rayleigh, Gamma, Poisson, and Bernoulli distri-

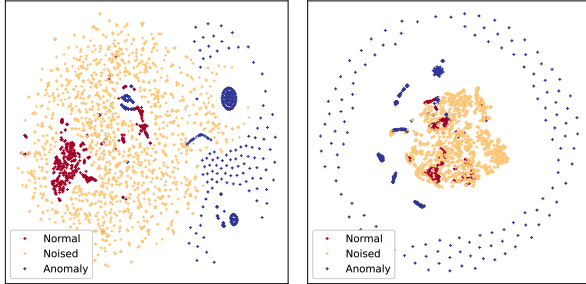
butions. For Salt&Pepper and Bernoulli noise, a probability vector from a uniform distribution generates a binary vector, dictating feature value alterations by changing values or reversing signs. Other distributions are adjusted to have zero mean and σ standard deviation. Detailed parameters are in described in Table 3, with results in Appendix H, Figure 10. Results show Salt&Pepper and Bernoulli noise performs poorly, indicating the effectiveness of our noise generation design. Gaussian, Rayleigh, and Uniform noise maintain stable performance with 92% AUC and 94% F1 scores, indicating our method’s robustness with various noise types.

Qualitative Visualization We utilize t-Distributed Stochastic Neighbor Embedding (t-SNE) visualization (Van der



(a) AUC and Its Rank on ResMLP with Different Noise Levels (b) F1 and Its Rank on ResMLP with Different Noise Levels (c) AUC and Its Rank on MLP with Different Noise Levels (d) F1 and Its Rank on MLP with Different Noise Levels

Figure 5: Sensitivity of Different Noise Level in $[0.1, 0.2, 0.5, 0.8, 1.0, 2.0, 3.0, 5.0]$. The mean rank (the lower the better) is calculated out of 8 noise levels.



(a) ResMLP T-SNE Visual. (b) MLP T-SNE Visualization

Figure 6: Qualitative visualization of the penultimate layer on the KDD-CUP99 dataset shows normal (red), noised (yellow), and anomalous (blue) instances. The noised instances help the model establish a decision boundary between in-distribution and out-of-distribution data. Real anomalies consistently fall outside this boundary in both ResMLP and MLP, confirming their effectiveness in anomaly detection.

Maaten and Hinton 2008) to show the allocation of normal, noised, and true anomalous samples in the neural network embedding space qualitatively. In Figure 6, we visualize the penultimate layer of the neural network using KDD-CUP99. It is seen that the introduction of noised instances into the dataset ostensibly aids the model in constructing a discriminative boundary that proficiently segregates the in-distribution data from its out-of-distribution counterparts. Empirical observations reveal that the actual anomalous data predominantly falls outside of this established boundary, a phenomenon consistently manifested in the experimental results for both ResMLP and MLP models.

Conclusion

In conclusion, we presented a novel noise evaluation-based method for unsupervised anomaly detection in tabular data. We assume the model can learn the anomalous pattern from the noised normal data. Predicting the magnitude of the noise shows inspiration on how much and where the abnormality is. We theoretically proved the generalizability and reliability of our method. Extensive experiments demonstrated that our approach outperforms other anomaly detection methods through 47 real tabular datasets in the UAD setting and 25

real tabular datasets in the OCC setting. An ablation study suggests that using Gaussian, Rayleigh, and Uniform noise has a stable performance. One potential limitation of this work is that we focused on tabular data only. Nevertheless, it is possible to extend our method to image data via the following two steps: 1) extract features of input images using a pretrained encoder; 2) apply our method to the extracted image features. This could be an interesting future work.

Acknowledgments

This work was supported by the Shenzhen Science and Technology Program under the Grant No.JCYJ20210324130208022 (Fundamental Algorithms of Natural Language Understanding for Chinese Medical Text Processing) and the General Program of Guangdong Basic and Applied Basic Research under Grant No.2024A1515011771.

References

- Abououf, M.; Singh, S.; Mizouni, R.; and Otkrok, H. 2023. Explainable AI for Event and Anomaly Detection and Classification in Healthcare Monitoring Systems. *IEEE Internet of Things Journal*.
- Aggarwal, C. C. 2016. *Outlier analysis second edition*.
- Angiulli, F.; and Pizzuti, C. 2002. Fast outlier detection in high dimensional spaces. In *European conference on principles of data mining and knowledge discovery*, 15–27. Springer.
- Anthony, M.; Bartlett, P. L.; Bartlett, P. L.; et al. 1999. *Neural network learning: Theoretical foundations*, volume 9. cambridge university press Cambridge.
- Bartlett, P. L.; Harvey, N.; Liaw, C.; and Mehrabian, A. 2019. Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research*, 20(1): 2285–2301.
- Behrmann, J.; Grathwohl, W.; Chen, R. T.; Duvenaud, D.; and Jacobsen, J.-H. 2019. Invertible residual networks. In *International conference on machine learning*, 573–582. PMLR.
- Ben-David, S.; Blitzer, J.; Crammer, K.; Kulesza, A.; Pereira, F.; and Vaughan, J. W. 2010. A theory of learning from different domains. *Machine learning*, 79: 151–175.

- Breunig, M. M.; Kriegel, H.-P.; Ng, R. T.; and Sander, J. 2000. LOF: identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, 93–104.
- Cai, C.; Hy, T. S.; Yu, R.; and Wang, Y. 2023. On the connection between mpnn and graph transformer. In *International Conference on Machine Learning*, 3408–3430. PMLR.
- Cai, J.; and Fan, J. 2022. Perturbation learning based anomaly detection. *Advances in Neural Information Processing Systems*, 35.
- Chang, C.-H.; Yoon, J.; Arik, S. Ö.; Udell, M.; and Pfister, T. 2023. Data-efficient and interpretable tabular anomaly detection. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 190–201.
- Chen, J.; Sathe, S.; Aggarwal, C.; and Turaga, D. 2017. Outlier detection with autoencoder ensembles. In *Proceedings of the 2017 SIAM international conference on data mining*, 90–98. SIAM.
- Chen, Y.; Tian, Y.; Pang, G.; and Carneiro, G. 2022. Deep one-class classification via interpolated gaussian descriptor. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 383–392.
- Ding, X.; Zhao, L.; and Akoglu, L. 2022. Hyperparameter sensitivity in deep outlier detection: Analysis and a scalable hyper-ensemble solution. *Advances in Neural Information Processing Systems*, 35: 9603–9616.
- Fan, J.; Chow, T. W. S.; and Qin, S. J. 2022. Kernel-Based Statistical Process Monitoring and Fault Detection in the Presence of Missing Data. *IEEE Transactions on Industrial Informatics*, 18(7): 4477–4487.
- Fu, D.; Zhang, Z.; and Fan, J. 2024. Dense projection for anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 8398–8408.
- Golan, I.; and El-Yaniv, R. 2018. Deep anomaly detection using geometric transformations. *Advances in neural information processing systems*, 31.
- Goyal, S.; Raghunathan, A.; Jain, M.; Simhadri, H. V.; and Jain, P. 2020. DROCC: Deep robust one-class classification. In *International conference on machine learning*, 3711–3721. PMLR.
- Han, S.; Hu, X.; Huang, H.; Jiang, M.; and Zhao, Y. 2022. Ad-bench: Anomaly detection benchmark. *Advances in Neural Information Processing Systems*, 35: 32142–32159.
- Hanneke, S.; et al. 2014. Theory of disagreement-based active learning. *Foundations and Trends® in Machine Learning*, 7(2-3): 131–309.
- He, Z.; Xu, X.; and Deng, S. 2003. Discovering cluster-based local outliers. *Pattern recognition letters*, 24(9-10): 1641–1650.
- Hendrycks, D.; Mazeika, M.; and Dietterich, T. 2018. Deep Anomaly Detection with Outlier Exposure. In *International Conference on Learning Representations*.
- Hendrycks, D.; Mazeika, M.; Kadavath, S.; and Song, D. 2019. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32.
- Hilal, W.; Gadsden, S. A.; and Yawney, J. 2022. Financial fraud: a review of anomaly detection techniques and recent advances. *Expert systems With applications*, 193: 116429.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Hu, T.; Zhang, J.; Yi, R.; Du, Y.; Chen, X.; Liu, L.; Wang, Y.; and Wang, C. 2024. Anomalydiffusion: Few-shot anomaly image generation with diffusion model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 8526–8534.
- Huang, Y.; Zhang, Y.; Wang, L.; Zhang, F.; and Lin, X. 2024. EntropyStop: Unsupervised Deep Outlier Detection with Loss Entropy. *ACM KDD*.
- Hussain, M.; Suh, J.-W.; Seo, B.-S.; and Hong, J.-E. 2023. How Reliable are the Deep Learning-based Anomaly Detectors? A Comprehensive Reliability Analysis of Autoencoder-based Anomaly Detectors. In *2023 Fourteenth International Conference on Ubiquitous and Future Networks (ICUFN)*, 317–322. IEEE.
- Jin, M.; Liu, Y.; Zheng, Y.; Chi, L.; Li, Y.-F.; and Pan, S. 2021. Anemone: Graph anomaly detection with multi-scale contrastive learning. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 3122–3126.
- Kelly, M.; Longjohn, R.; and Nottingham, K. ????. The UCI Machine Learning Repository. Accessed: 2023-12-27.
- Kifer, D.; Ben-David, S.; and Gehrke, J. 2004. Detecting change in data streams. In *VLDB*, volume 4, 180–191. Toronto, Canada.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Li, Z.; Zhao, Y.; Botta, N.; Ionescu, C.; and Hu, X. 2020. COPOD: copula-based outlier detection. In *2020 IEEE international conference on data mining (ICDM)*, 1118–1123. IEEE.
- Li, Z.; Zhao, Y.; Hu, X.; Botta, N.; Ionescu, C.; and Chen, G. H. 2023. ECOD: Unsupervised Outlier Detection Using Empirical Cumulative Distribution Functions. *IEEE Transactions on Knowledge and Data Engineering*, 35(12): 12181–12193.
- Liu, F. T.; Ting, K. M.; and Zhou, Z.-H. 2008. Isolation forest. In *2008 eighth IEEE international conference on data mining*, 413–422. IEEE.
- Liu, Y.; Li, Z.; Zhou, C.; Jiang, Y.; Sun, J.; Wang, M.; and He, X. 2019. Generative adversarial active learning for unsupervised outlier detection. *IEEE Transactions on Knowledge and Data Engineering*, 32(8): 1517–1528.
- Liznerski, P.; Ruff, L.; Vandermeulen, R. A.; Franks, B. J.; Kloft, M.; and Muller, K. R. 2020. Explainable Deep One-Class Classification. In *International Conference on Learning Representations*.
- Pang, G.; Shen, C.; Cao, L.; and Hengel, A. V. D. 2021. Deep learning for anomaly detection: A review. *ACM computing surveys (CSUR)*, 54(2): 1–38.

- Papadimitriou, S.; Kitagawa, H.; Gibbons, P. B.; and Faloutsos, C. 2003. Loci: Fast outlier detection using the local correlation integral. In *Proceedings 19th international conference on data engineering (Cat. No. 03CH37405)*, 315–326. IEEE.
- Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic differentiation in pytorch.
- Qiu, C.; Pfrommer, T.; Kloft, M.; Mandt, S.; and Rudolph, M. 2021. Neural transformation learning for deep anomaly detection beyond images. In *International conference on machine learning*, 8703–8714. PMLR.
- Reddi, S. J.; Kale, S.; and Kumar, S. 2018. On the Convergence of Adam and Beyond. In *International Conference on Learning Representations*.
- Roth, K.; Pemula, L.; Zepeda, J.; Schölkopf, B.; Brox, T.; and Gehler, P. 2022. Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14318–14328.
- Ruff, L.; Vandermeulen, R.; Goernitz, N.; Deecke, L.; Siddiqui, S. A.; Binder, A.; Müller, E.; and Kloft, M. 2018. Deep one-class classification. In *International conference on machine learning*, 4393–4402. PMLR.
- Saeed, M. M.; Saeed, R. A.; Abdelhaq, M.; Alsaqour, R.; Hasan, M. K.; and Mokhtar, R. A. 2023. Anomaly Detection in 6G Networks Using Machine Learning Methods. *Electronics*, 12(15): 3300.
- Schlegl, T.; Seeböck, P.; Waldstein, S. M.; Schmidt-Erfurth, U.; and Langs, G. 2017. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, 146–157. Springer.
- Schölkopf, B.; Platt, J. C.; Shawe-Taylor, J.; Smola, A. J.; and Williamson, R. C. 2001. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7): 1443–1471.
- Shenkar, T.; and Wolf, L. 2022. Anomaly detection for tabular data with internal contrastive learning. In *International Conference on Learning Representations*.
- Siddiqui, M. A.; Stokes, J. W.; Seifert, C.; Argyle, E.; McCann, R.; Neil, J.; and Carroll, J. 2019. Detecting cyber attacks using anomaly detection with explanations and expert feedback. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2872–2876. IEEE.
- Sohn, K.; Li, C.-L.; Yoon, J.; Jin, M.; and Pfister, T. 2020. Learning and Evaluating Representations for Deep One-Class Classification. In *International Conference on Learning Representations*.
- Song, Y.; and Ermon, S. 2019. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.
- Song, Y.; and Ermon, S. 2020. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33: 12438–12448.
- Stolfo, S.; Fan, W.; Lee, W.; Prodromidis, A.; and Chan, P. 1999. KDD Cup 1999 Data. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C51C7N>.
- Tax, D. M.; and Duin, R. P. 2004. Support vector data description. *Machine learning*, 54: 45–66.
- Touvron, H.; Bojanowski, P.; Caron, M.; Cord, M.; El-Nouby, A.; Grave, E.; Izacard, G.; Joulin, A.; Synnaeve, G.; Verbeek, J.; et al. 2022. Resmlp: Feedforward networks for image classification with data-efficient training. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4): 5314–5321.
- Tur, A. O.; Dall’Asen, N.; Beyan, C.; and Ricci, E. 2023. Exploring diffusion models for unsupervised video anomaly detection. In *2023 IEEE International Conference on Image Processing (ICIP)*, 2540–2544. IEEE.
- Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).
- Vincent, P. 2011. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7): 1661–1674.
- Wang, S.; Wang, X.; Zhang, L.; and Zhong, Y. 2021. Auto-AD: Autonomous hyperspectral anomaly detection network based on fully convolutional autoencoder. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1–14.
- Xiao, F.; Zhou, J.; Han, K.; Hu, H.; and Fan, J. 2025. Unsupervised anomaly detection using inverse generative adversarial networks. *Information Sciences*, 689: 121435.
- Yan, X.; Zhang, H.; Xu, X.; Hu, X.; and Heng, P.-A. 2021. Learning semantic context from normal samples for unsupervised anomaly detection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, 3110–3118.
- Yang, X.; Latecki, L. J.; and Pokrajac, D. 2009. Outlier detection with globally optimal exemplar-based GMM. In *Proceedings of the 2009 SIAM international conference on data mining*, 145–154. SIAM.
- Yang, X.; Qi, X.; and Zhou, X. 2023. Deep Learning Technologies for Time Series Abnormality Detection in Healthcare: A Review. *IEEE Access*.
- Zavrtanik, V.; Kristan, M.; and Skočaj, D. 2021. Draem-a discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, 8330–8339.
- Zhang, X.; Li, N.; Li, J.; Dai, T.; Jiang, Y.; and Xia, S.-T. 2023. Unsupervised surface anomaly detection with diffusion probabilistic model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6782–6791.
- Zhang, Y.; Sun, Y.; Cai, J.; and Fan, J. 2024. Deep Orthogonal Hypersphere Compression for Anomaly Detection. In *Proceedings of the International Conference on Learning Representations*.
- Zhao, Y.; Nasrullah, Z.; and Li, Z. 2019. PyOD: A Python Toolbox for Scalable Outlier Detection. *Journal of Machine Learning Research*, 20(96): 1–7.
- Zong, B.; Song, Q.; Min, M. R.; Cheng, W.; Lumezanu, C.; Cho, D.; and Chen, H. 2018. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International conference on learning representations*.

Supplementary Material for: Unsupervised Anomaly Detection for Tabular Data Using Noise Evaluation

Appendix A: Connection Between Noise Evaluation and Denoising Score Matching

We first review the denoising score matching (Song and Ermon 2019; Vincent 2011). The *score* of a probability density $p(\mathbf{x})$ is defined as $\nabla_{\mathbf{x}} \log p(\mathbf{x})$. Score matching train a score network parametrized by θ , $s_{\theta}(\mathbf{x})$, to estimate $\nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})$. Denoising score matching is a variant of score matching. It first perturbs the data \mathbf{x} with a known noise distribution $q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})$ and estimate the score of the perturbed data distribution $q_{\sigma}(\hat{\mathbf{x}}) \triangleq \int q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})p_{data}(\mathbf{x})d\mathbf{x}$. The learning objective was proved equivalent to the following:

$$\min_{\theta} \frac{1}{2} \mathbb{E}_{q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})} \|s_{\theta}(\hat{\mathbf{x}}) - \nabla_{\hat{\mathbf{x}}} \log q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})\|_2^2. \quad (12)$$

Here we claim optimizing the noise evaluation learning objective is a lower bound of the denoising score matching.

Proof. Without loss of generality, similar to (Vincent 2011), we consider the perturbation kernel $q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})$ satisfies

$$\frac{\partial q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})}{\partial \hat{\mathbf{x}}} = \Delta(\mathbf{x} - \hat{\mathbf{x}}), \Delta > 0.$$

If the perturbation kernel is Gaussian noise, $\Delta = \frac{1}{\sigma^2}$.

Let us now choose the model p as

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z(\theta)} \exp(-f_{\theta}(\mathbf{x})),$$

$$f_{\theta}(\mathbf{x}) = -\Delta \left(\int h_{\theta}(\mathbf{x}) d\mathbf{x} \right), h_{\theta}(\mathbf{x}) \geq \mathbf{0},$$

where $Z(\theta)$ is an intractable partition function and h is our noise evaluation model. Then,

$$s_{\theta}(\mathbf{x}) = \frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \mathbf{x}} = -\frac{f_{\theta}(\mathbf{x})}{\partial \mathbf{x}} = \Delta h_{\theta}(\mathbf{x}).$$

Since the noised sample in our method is generated as $\hat{\mathbf{x}} = \mathbf{x} + \epsilon$, then we have.

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})} \|s_{\theta}(\hat{\mathbf{x}}) - \nabla_{\hat{\mathbf{x}}} \log q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})\|_2^2 \\ &= \frac{1}{2} \Delta \mathbb{E}_{q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})} \|h_{\theta}(\hat{\mathbf{x}}) - (\mathbf{x} - \hat{\mathbf{x}})\|_2^2 \\ &\geq \frac{1}{2} \Delta \mathbb{E}_{q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})} \|h_{\theta}(\hat{\mathbf{x}}) - |\mathbf{x} - \hat{\mathbf{x}}|\|_2^2 \\ &= \frac{1}{2} \Delta \mathbb{E}_{q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x})} \|h_{\theta}(\hat{\mathbf{x}}) - |\epsilon|\|_2^2 \end{aligned}$$

For the input data from the normal class, the noise magnitude is naturally $\mathbf{0}$. According to (Song and Ermon 2019), setting $\epsilon = \mathbf{0}$ in the learning objective makes the model can learn the underlying data distribution, i.e., $s_{\theta}(\mathbf{x}) = \nabla_{\mathbf{x}} \log q_{\sigma}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})$. Hence, we show that the noise evaluation learning objective is a lower bound of the denoising score matching. Q.E.D.

Since the goal of our method is anomaly detection instead of sample generation, predicting the magnitude of the noise is enough to detect samples far from the $p_{data}(\mathbf{x})$ which are usually anomalies. For samples belonging to the normal class, it is natural that the gradient of its probability density is 0. We provide an intuitive example here with Gaussian distribution. Suppose $x \sim \mathcal{N}(0, \sigma)$, its gradient of log density function $\frac{d}{dx} \log p(x) = -\frac{1}{\sigma^2}x$. We see if the magnitude of the gradient is large, the sample is far away from the distribution center 0. In our method, we involve diverse noise with different variances because $p_{data}(\mathbf{x})$ is unknown and $\frac{\partial q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x}_j)}{\partial \hat{\mathbf{x}}}$ and $\frac{\partial q_{\sigma}(\hat{\mathbf{x}}|\mathbf{x}_i)}{\partial \hat{\mathbf{x}}}$ maybe very different for the same $\hat{\mathbf{x}}$ with different i, j . Since we want to enlarge such gradients, we use various noise levels. It is similar to (Song and Ermon 2019, 2020) using an annealed noise level. Therefore, noise evaluation can detect anomaly data from the perspective of score matching.

Appendix B: Noise Evaluation Training Algorithm

Algorithm 2: Noise Evaluation Training

Input: Training dataset \mathcal{X} , max. iterations T ;

- 1: Normalize \mathcal{X} to have mean of 0 and standard deviation of 1;
- 2: Initialize model θ ;
- 3: $t \leftarrow 1$;
- 4: **while** $t \leq T$ **do**
- 5: **for** each batch of data $\mathcal{B} \in \mathcal{X}$ **do**
- 6: //get batch size and dimension.
- 7: $b, d \leftarrow \mathcal{B}.shape$;
- 8: //Algorithm 1.
- 9: $\mathcal{E} \leftarrow \text{NoiseGeneration}(b, d)$;
- 10: $\hat{\mathcal{B}} \leftarrow \mathcal{B} + \mathcal{E}$;
- 11: $l \leftarrow \frac{\sum_i \|h_{\theta}(\mathbf{x}_i)\|_2^2}{\sum_i \|h_{\theta}(\hat{\mathbf{x}}_i) - |\epsilon_i|\|_2^2}, \mathbf{x}_i \in \mathcal{B}, \hat{\mathbf{x}}_i \in \hat{\mathcal{B}}, \epsilon_i \in \mathcal{E}$;
 //Equation (9).
- 12: Optimize θ by Adam (Kingma and Ba 2014; Reddi, Kale, and Kumar 2018) optimizer;
- 13: **end for**
- 14: $t \leftarrow t + 1$;
- 15: **end while**

Output: Optimized model h_{θ} ;

The overall training algorithm for the noise evaluation training is in Algorithm 2. Notice that for each training epoch, we randomly generate a new noised instance for each training sample. This helps us enlarge the sampling number from the noise distribution for better learning. In line 9, there are some optional noise generation schemes such as using different noise types, different noise levels, and different noise ratios.

Time and Space Complexity Let b denote the batch size, w_{\max} represent the maximum width of the hidden layers in an L -layer neural network, and d indicate the dimension of the input data. Under these parameters, the time complexity of the proposed methods is at most $\mathcal{O}(bdw_{\max}LT)$, where T is the maximum number of iterations. The space complexity

is at most $\mathcal{O}(bd + dw_{\max} + (L-1)w_{\max}^2)$. These complexities scale linearly with the number of samples, which underscores the scalability of the proposed methods to large datasets. Moreover, in the case of high-dimensional data (i.e., large d), choosing a smaller w_{\max} can further enhance computational efficiency.

Appendix C: Proof of Theoretical Analysis

Proof for Proposition 1

Proof. For convenience, denote the noise variable as E , then $\hat{X} = X + E$. We have

$$\begin{aligned} & H(\hat{X} | X) \\ &= - \sum_x p(X=x) \sum_{\hat{x}} p(\hat{X}=\hat{x} | X=x) \log p(\hat{X}=\hat{x} | X=x) \\ &= - \sum_x p(X=x) \\ & \quad \sum_{\hat{x}} p(E=\hat{x}-x | X=x) \log p(E=\hat{x}-x | X=x) \\ &= - \sum_x p(X=x) \sum_{e'} p(E=e' | X=x) \log p(E=e' | X=x) \\ &= H(E | X) \end{aligned} \tag{13}$$

Since E and X are independent, we have $H(E | X) = H(E)$. On the other hand, we have $H(\hat{X}) \geq H(\hat{X}|X)$, which can be proved by the definition of conditional entropy and Jensen's inequality (or log sum inequality more specifically). Then we arrive at

$$H(\hat{X}) \geq H(\hat{X} | X) = H(E | X) = H(E).$$

Similarly, by symmetry, we have

$$H(\hat{X}) \geq H(\hat{X} | E) = H(X | E) = H(X).$$

This inequality holds for each variable in \mathbf{x} . Then the total entropy of the d variables increases, which means the data becomes more disordered. Q.E.D.

Definition of \mathcal{H} -divergence

Definition 2. (Based on Kifer et al., 2004 (Kifer, Ben-David, and Gehrke 2004)) Given a domain \mathcal{X} with D and D' probability distributions over \mathcal{X} , let \mathcal{H} be a hypothesis class on \mathcal{X} and denote by $I(h)$ the set for which $h \in \mathcal{H}$ is the characteristic function; that is, $x \in I(h) \iff h(x) = 1$. The \mathcal{H} -divergence between D and D' is

$$d_{\mathcal{H}}(D, D') = 2 \sup_{h \in \mathcal{H}} \left| \Pr_D[I(h)] - \Pr_{D'}[I(h)] \right|.$$

The definition is defined in (Ben-David et al. 2010). We repeat here to provide reader convenience. In addition, $\mathcal{H}\Delta\mathcal{H}$ is a symmetric difference hypothesis space, where $f \in \mathcal{H}\Delta\mathcal{H}$, $f(\mathbf{x}) = I(g(h(\mathbf{x})) > \tau) \text{ XOR } I(g(h'(\mathbf{x})) > \tau)$ for some $h, h' \in \mathcal{H}$.

Proof of Theorem 1

Before proving the theorem, we provide the following necessary lemma according to (Ben-David et al. 2010; Bartlett et al. 2019).

Lemma 1. Let \mathcal{H} be a hypothesis space on \mathbb{R}^d with VC dimension d_{vc} . $h \in \mathcal{H}$ is a ReLU-activated deep neural network with L layers, and the number of neurons at each layer is in the order of p . If \mathcal{X} and \mathcal{X}' are samples of size N from two distributions \mathcal{D} and \mathcal{D}' respectively and $\hat{d}_{\mathcal{H}}(\mathcal{X}, \mathcal{X}')$ is the empirical \mathcal{H} -divergence between samples, then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$d_{\mathcal{H}}(\mathcal{D}, \mathcal{D}') \leq \hat{d}_{\mathcal{H}}(\mathcal{X}, \mathcal{X}') + 4 \sqrt{\frac{d_{vc} \log(2N) + \log(\frac{2}{\delta})}{N}},$$

where $d_{vc} = \mathcal{O}(pL \log(pL))$.

Proof. We slightly modify the lemma in (Ben-David et al. 2010) and Theorem 3.4 in (Kifer, Ben-David, and Gehrke 2004) with explicit definition of our neural network and specified VC dimension $\mathcal{O}(pL \log(pL))$. Q.E.D.

Following (Ben-David et al. 2010), we have an additional lemma.

Lemma 2. For any hypotheses $h, h' \in \mathcal{H}$,

$$|\varepsilon(h, h') - \varepsilon(h, h')| \leq \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}, \mathcal{D}').$$

Proof. By the definition of $\mathcal{H}\Delta\mathcal{H}$ -distance,

$$\begin{aligned} & d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}, \mathcal{D}') \\ &= 2 \sup_{h, h' \in \mathcal{H}} \left| \Pr_{x \sim \mathcal{D}} [h(x) \neq h'(x)] - \Pr_{x \sim \mathcal{D}'} [h(x) \neq h'(x)] \right| \\ &= 2 \sup_{h, h' \in \mathcal{H}} |\varepsilon(h, h') - \varepsilon'(h, h')| \geq 2 |\varepsilon(h, h') - \varepsilon'(h, h')|. \end{aligned}$$

Q.E.D.

Then, we have the proof.

Proof. Define

$$\varepsilon_{\mathcal{D}}(h, h^*) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [|I(g(h_{\theta}(\mathbf{x})) > \tau) - I(g(h_{\theta}^*(\mathbf{x})) > \tau)|].$$

We have

$$\begin{aligned} \varepsilon_{\tilde{\mathcal{D}}_H}(h) &\leq \varepsilon_{\tilde{\mathcal{D}}_H}(h^*) + \varepsilon_{\tilde{\mathcal{D}}_H}(h, h^*) \\ &\leq \varepsilon_{\tilde{\mathcal{D}}_H}(h^*) + \varepsilon_{\tilde{\mathcal{D}}}(h, h^*) + |\varepsilon_{\tilde{\mathcal{D}}_H}(h, h^*) - \varepsilon_{\tilde{\mathcal{D}}}(h, h^*)| \\ &\leq \varepsilon_{\tilde{\mathcal{D}}_H}(h^*) + \varepsilon_{\tilde{\mathcal{D}}}(h, h^*) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\hat{\mathcal{D}}, \tilde{\mathcal{D}}_H) \\ &\leq \varepsilon_{\tilde{\mathcal{D}}_H}(h^*) + \varepsilon_{\tilde{\mathcal{D}}}(h) + \varepsilon_{\tilde{\mathcal{D}}}(h^*) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\hat{\mathcal{D}}, \tilde{\mathcal{D}}_H) \\ &= \varepsilon_{\tilde{\mathcal{D}}}(h) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\hat{\mathcal{D}}, \tilde{\mathcal{D}}_H) + \lambda \\ &\leq \varepsilon_{\tilde{\mathcal{D}}}(h) + \frac{1}{2} \hat{d}_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{X}}_H, \hat{\mathcal{X}}) + 4 \sqrt{\frac{2d_{vc} \log(2N) + \log(\frac{2}{\delta})}{N}} + \lambda \end{aligned}$$

The VD dimension of $\mathcal{H}\Delta\mathcal{H}$ is at most twice the VC dimension of \mathcal{H} (Anthony et al. 1999). Q.E.D.

Proof of Theorem 2

Proof. Assumption 1 indicates that

$$c\|\mathbf{x} - \tilde{\mathbf{x}}\| \leq |g(h_{\theta}(\mathbf{x})) - g(h_{\theta}(\tilde{\mathbf{x}}))|.$$

According to the definition of d_{\min} , we have

$$g(h_{\theta}(\tilde{\mathbf{x}})) \geq cd_{\min} + g(h_{\theta}(\mathbf{x})) \quad (14)$$

or

$$g(h_{\theta}(\tilde{\mathbf{x}})) \leq g(h_{\theta}(\mathbf{x})) - cd_{\min}. \quad (15)$$

Since $\epsilon < cd_{\min}$, we have $g(h_{\theta}(\tilde{\mathbf{x}})) < 0$, which contradicts the fact that $g(h_{\theta}(\tilde{\mathbf{x}}))$ is always positive. Therefore, (15) will not happen and we will only have (14). It follows that

$$g(h_{\theta}(\tilde{\mathbf{x}})) \geq cd_{\min}. \quad (16)$$

This means if $cd_{\min} > \tau$, $\tilde{\mathbf{x}}$ can be detected successfully. Putting the conditions $\epsilon < cd_{\min}$ and $cd_{\min} > \tau$, $\tilde{\mathbf{x}}$ together, we complete the proof. Q.E.D.

Dataset	# Sample	Dims.	# Class	% Anomaly
abalone	4177	8	28	94.19
arrhythmia	452	280	2	45.80
breastw	699	10	2	34.48
cardio	2126	22	2	22.15
ecoli	336	8	8	77.16
glass	214	10	6	72.74
ionosphere	351	35	2	35.90
kdd	5209460	39	2	80.52
letter	20000	17	26	96.15
lympho	142	19	2	42.96
mammography	11183	7	2	2.32
mulcross	262144	5	2	10.00
musk	7074	167	2	17.30
optdigits	5620	65	10	90.00
pendigits	10992	17	10	90.00
pima	768	9	2	34.90
satimage	6435	37	6	83.33
seismic	2584	15	2	6.58
shuttle	58000	10	7	75.03
speech	3686	401	2	1.65
thyroid	7200	7	2	7.29
vertebral	310	7	2	67.74
vowels	1456	13	2	3.43
wbc	569	31	2	37.26
wine	178	14	3	66.67

Table 4: 25 real-world tabular datasets tested in this paper under OCC setting. For datasets with more than 2 labeled classes, an average anomaly ratio is recorded. Kdd refers to KDD-CUP99 (Stolfo et al. 1999).

Appendix D: Dataset Summary

We conducted our experiment under the UAD dataset setting with 47 benchmark datasets commonly used in UAD research. The dataset information is shown in Table 5. As for the OCC dataset setting. The detailed dataset information including the

number of samples, dimensionality, the number of classes, and the percent of anomaly is summarized in Table 4. Note that there are datasets with the same name. However, due to the different dataset processing and training data splitting. They are completely different under OCC setting from the UAD setting.

Dataset	# Sample	Dim.	% Anomaly
ALOI	49534	27	3.04
anthyroid	7200	6	7.42
backdoor	95329	196	2.44
breast	683	9	34.99
campaign	41188	62	11.27
cardio	1831	21	9.61
Cardiotocography	2114	21	22.04
celeba	202599	39	2.24
census	299285	500	0.60
cover	286048	10	6.20
donors	619326	10	5.93
fault	1941	27	34.67
fraud	284807	29	0.17
glass	214	9	42.21
Hepatitis	80	17	16.25
http	567498	3	0.39
InternetAds	1966	1555	18.72
Ionosphere	351	32	35.90
landsat	6435	36	20.71
letter	1600	32	15.10
Lymphography	148	18	4.05
magic	19020	10	35.16
mammography	11183	6	2.32
mnist	7603	100	9.21
musk	3062	166	3.17
optdigits	5216	64	2.88
PageBlocks	5393	10	9.46
pendigits	6870	16	2.27
Pima	768	8	34.90
satellite	6435	36	31.64
satimage-2	5803	36	1.22
shuttle	49097	9	7.15
skin	245057	3	20.75
smtp	95156	3	0.03
SpamBase	4207	57	39.91
speech	3686	400	1.65
Stamps	340	9	9.12
thyroid	3772	6	2.47
vertebral	240	6	12.50
vowels	1456	12	3.43
Waveform	3443	21	2.90
WBC	223	9	4.48
WDBC	367	30	2.72
Wilt	4819	5	5.33
wine	129	13	7.75
WPBC	198	33	23.74
yeast	1484	8	34.16

Table 5: 47 real-world tabular datasets tested under UAD setting.

Appendix E: Neural Network Architecture

VanillaMLP

The VanillaMLP is a plain Relu-activated feed-forward neural network with 4 fully connected layers. The architecture is shown in Figure 7a. The input dimension is d . If $d \leq 64$, the size of the hidden layer is set to 64. If $d > 64$, the size of the hidden layer is set to 256. In the experiments, datasets, arrhythmia, musk, and speech, are with $d > 64$.

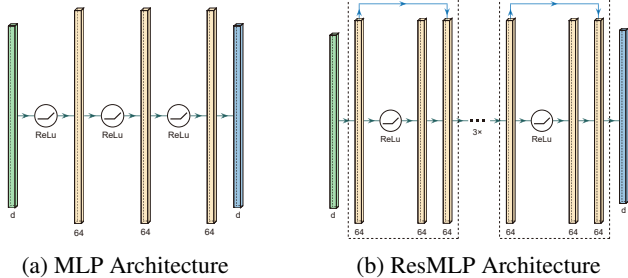


Figure 7: Deep Neural Network Architectures.

ResMLP

The ResMLP is a Relu-activated feed-forward neural network with 5 residual blocks. The architecture is shown in Figure 7b. The input dimension is d . If $d \leq 64$, the size of the hidden layer is set to 64. If $d > 64$, the size of the hidden layer is set to 256. In the experiments, datasets, arrhythmia, musk, and speech, are with $d > 64$.

Appendix F: Detailed Anomaly Detection Results

The detailed results in terms of AUC score and F1 score under the UAD setting are reported here in Table 6 and 7. We do not include the results for AnoGAN baseline because it is not efficient to run. So, we only compare 11 baselines for the UAD setting. The paired t-test is conducted based on the average AUC/F1 score through 10 runs in Tables 2, 5, 6, and 7. For example, in Table 2, every AD method has 25 values (as there are 25 datasets). Thus for any two AD methods, we have 25 pairs of values, on which we performed the paired t-test. The last two lines in the table are the p-values of our two methods (ResMLP and MLP) against each baseline method. Since the anomaly in the test set is rare, F1 score metric is more important to show the significance of our method. It is seen that our method has the highest mean F1 score (69.40%) and average rank (2.70). Compared with other baseline methods, our method has statistically significant improvement in terms of pair t-test p-value.

Appendix G: F1 Score Results on OCC setting

The performance in terms of F1 score under the OCC setting is presented in Table 8. Our approach also obtains the highest mean F1 score and mean rank out of 12 baselines.

Appendix H: Ablation Experiment Results

Sensitivity of Different Noise Levels

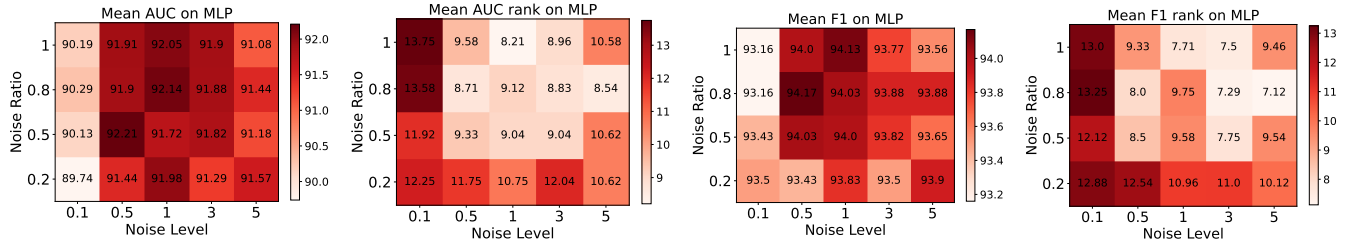
To explore how different noise levels contribute to the performance, we keep the noise type as Gaussian noise. In each training batch, we also generate 3 noised instances which is consistent with the former setting. We search through different noise levels in $[0.1, 0.2, 0.5, 0.8, 1.0, 2.0, 3.0, 5.0]$. We report the average AUC, F1 score, and rank (out of 8 different noise levels) in Figure 5. We observe that in all results if the noise level is too small, the average performance is low. This indicates noise with a small level would confuse the model because the distance between the normal samples and anomalies is very small. On the other hand, excessively high noise levels result in an expanded output value range, thereby enlarging the sampling space. This enlargement impacts the expressiveness of the model, reducing its effectiveness. However, when the noise level grows too large the score is also decreased. This means a large noise level has a large learning target value and a large sampling space, which may result in under-fitting training. The best choice of noise level appears around 1.0. From the perspective of denoising score matching(Song and Ermon 2019), a higher noise level will lead to inaccurate score estimation while a lower noise level will make it hard to learn data distribution in a relatively low-density region. Hence, a diverse noise level selection is preferred.

Sensitivity of Different Noise Ratios

The noise ratio means the percentage of noise appearing the feature in a sample. For different noise ratios, we keep utilizing Gaussian noise and the same noise level generation as the default. In each training batch, we generate 1 noised instances with the specific noise ratio instead. We vary the noise ratio in $[0.2, 0.5, 0.8, 1.0]$. We plot a heatmap jointly with different noise levels in Figure and 8 and 9. We observe an augmentation in model accuracy concomitant with an increase in the noise ratio across the spectrum of noise levels tested. This trend suggests that our model exhibits an affinity for a higher noise ratio in its training regimen, implying that a moderated presence of noise may be beneficial to the model’s training efficacy and predictive accuracy. In ResMLP result, similar results are observed that a higher noise ratio is beneficial to the model performance. Noise levels either too high or too low can lead to a bad performance.

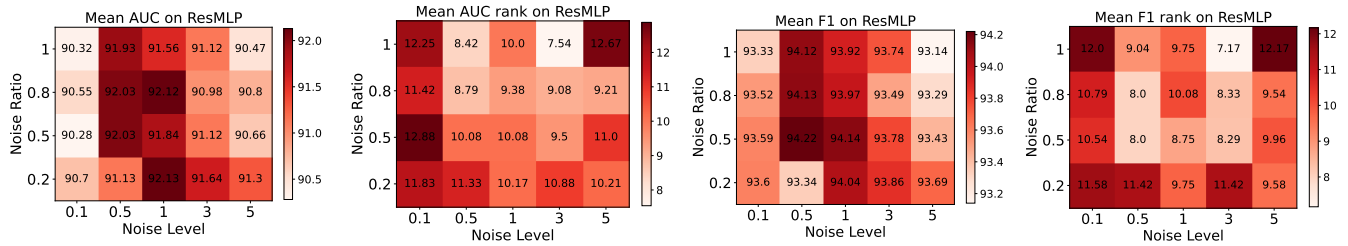
Sensitivity of Different Noise Type

We study utilizing different noise types to generate the noise. For a fair comparison, we sample all noise with 0 mean and σ standard deviation (noise level). We adopt several common noise distributions, including Salt&Pepper, Gaussian, Laplace, Uniform, Rayleigh, Gamma, Poisson, and Bernoulli. For incorporating noise of a Bernoulli nature, the initial step involves the generation of a probability vector, which is derived from a uniform distribution. This vector is then utilized to produce a corresponding binary vector. The binary vector dictates the positions at which the original feature values will be altered—should the noise be present, a reversal of the feature’s sign is executed. For Salt&Pepper, similar to Bernoulli,



(a) AUC on MLP with Different Noise Levels and Ratios (b) Rank by AUC on MLP with Different Noise Levels and Ratios (c) F1 on MLP with Different Noise Levels and Ratios (d) Rank by F1 on MLP with Different Noise Levels and Ratios

Figure 8: Sensitivity of Different Noise Levels and Ratio on MLP. We test different noise levels in $[0.1, 0.5, 1.0, 3.0, 5.0]$, and noise ratios in $[0.2, 0.8, 0.5, 1.0]$. The mean rank (the lower the better) is calculated out of 20 different settings.



(a) AUC on ResMLP with Different Noise Levels and Ratios (b) Rank by AUC on ResMLP with Different Noise Levels and Ratios (c) F1 on ResMLP with Different Noise Levels and Ratios (d) Rank by F1 on ResMLP with Different Noise Levels and Ratios

Figure 9: Sensitivity of Different Noise Levels and Ratio on ResMLP. We test different noise levels in $[0.1, 0.5, 1.0, 3.0, 5.0]$, and noise ratios in $[0.2, 0.8, 0.5, 1.0]$. The mean rank (the lower the better) is calculated out of 20 different settings.

we replace the value with the maximum or minimum value randomly in a batch instead of flipping the sign. For the other distributions, we adjust the parameter to make the generated noise having zero mean and σ standard deviation. The detailed generation parameter is shown in Table 3. The result is illustrated in Figure 10. Note that gamma1 and gamma3 represent gamma distribution with hyper-parameter $b = 1$, and 3, respectively. In the results, we observe that Bernoulli and Salt&Pepper noise has a low performance score. It indicates the effectiveness of involving different noise levels in our noise generation design. Gaussian, Rayleigh, and Uniform noise have relatively stable performance with 92% AUC score and 94% F1 score. It reveals that our method is robust with multiple types of noise.

Methods	Time (s)
Ours	5.32 ± 0.014
PLAD	6.22 ± 0.382

Table 9: Comparative study on the time cost for noise generation. Our method takes less time to generate the noised sample in an epoch (3800 batches). In addition, our method can be implemented in a pre-generation manner to further save time.

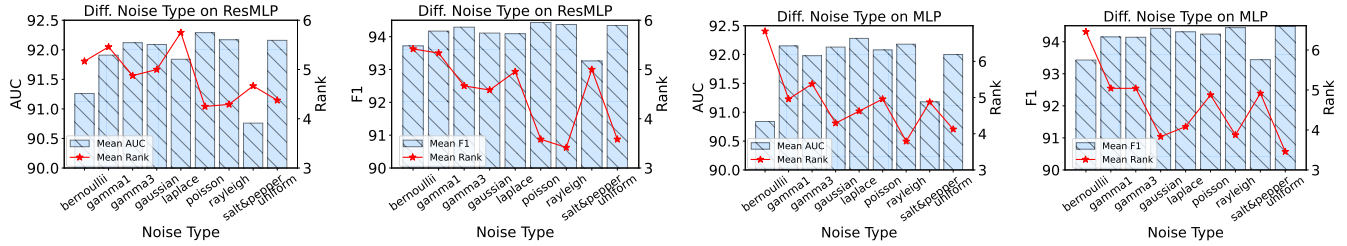
Methods	Training Time (s)	Inference Time (s)
Ours	37.19 ± 1.33	6.41 ± 0.27
PLAD (2022)	40.54 ± 0.38	4.33 ± 0.08
NeuTraLAD (2021)	39.72 ± 1.14	6.79 ± 0.24
DPAD (2024)	26.84 ± 0.24	287.6 ± 25.1
SCAD (2022)	37.17 ± 0.61	6.46 ± 0.28

Table 10: Comparative study of the time cost for model training and inference on KDDCUP-99 dataset. All methods use the same 5-layer backbone MLP. Our method takes less time to train than PLAD and NeuTraLAD. DPAD takes the least time for training because it is a density estimation method, but it takes unacceptable time for inference. Note that our method can reach a faster training speed by utilizing the pre-generated noisy sample.

Appendix I: Time Cost Analysis

Time Cost for Noise Generation

Suppose the batch size is b and the data dimensionality is d , the noise generation time complexity is $\mathcal{O}(bd)$ according to Algorithm 1. In contrast, other methods involving perturbation (Cai and Fan 2022; Qiu et al. 2021) and adversarial sample (Goyal et al. 2020) have time complexity with $\mathcal{O}(bdW)$, where W is workload related to a neural network module. Hence, our method for processing the generated negative sample is more efficient. We report the average time cost



(a) AUC and Rank on ResMLP with (b) F1 and Rank on ResMLP with (c) AUC and Rank on MLP with (d) F1 and Rank on MLP with Different Noise Type Different Noise Type Different Noise Type Different Noise Type

Figure 10: Sensitivity of Different Noise Types. We test different noise types including Salt&Pepper, Gaussian, Laplace, Uniform, Rayleigh, Gamma, Poisson, and Bernoulli distributions. The mean rank (the lower the better) is calculated out of 9 different noise types. The results show that Gaussian, Rayleigh, and Uniform distributions have relatively stable performance.

for generating the noise sample per epoch on KDDCUP-99 dataset in Table 9. The time cost for our noise generation is 5.69 ± 0.014 seconds while the negative sample processing time in PLAD is 6.22 ± 0.38 seconds. In addition, our noised sample can be generated before the training started. Hence, the generation time can be negligible.

Time Cost for Training and Inference

We also compare the training and inference time with other recent deep learning based methods, PLAD, NeuTraLAD, DPAD, and SCAD, on KDDCUP-99 dataset. The results are reported in Table 10. All methods use the same 5-layer backbone MLP. Our method takes less time to train than PLAD and NeuTraLAD. DPAD takes the least time for training because it is a density estimation method, but it takes unacceptable time for inference. Note that our method can reach a faster training speed by utilizing the pre-generated noisy sample. Hence, our method is more efficient at the training stage.

Methods	Hyper-param.	#
Ours	Noise Level, Ratio, and Type	3
PLAD (Cai and Fan 2022)	λ , Types of perturbator	2
NeuTraLAD (Qiu et al. 2021)	τ, K, m	3
DPAD (Fu, Zhang, and Fan 2024)	m, γ, λ, k	4
SCAD (Shenkar and Wolf 2022)	k, u, τ, r	4
DROCC (Goyal et al. 2020)	r, λ, μ, η	4
DOHSC (Zhang et al. 2024)	k, k', λ, μ, ν	5

Table 11: The comparison of the number of hyper-parameters. Although PLAD has fewer hyper-parameters, the choice of different types of perturbator would lead to a completely different implementation and would result in more new hyper-parameters.

Appendix J: Number of Hyper-parameters

We compare the number of hyper-parameters with other recent deep learning based UAD methods. The comparison is listed in Table 11. Note that all hyper-parameters of our method are solely related to noise generation, and no

hyper-parameter is involved in the learning objective except for the network training design (which exists in all deep learning based methods). Although PLAD has fewer hyper-parameters, the choice of different types of perturbator would lead to a completely different implementation and would result in more new hyper-parameters. Hence, our method is easy to implement.

Discussions

Interpretability For those deep learning-based methods (Golan and El-Yaniv 2018; Hendrycks et al. 2019; Hendrycks, Mazeika, and Dietterich 2018) that only output an anomaly score, it is hard to explain the forwarding non-linear process of the neural network black box. In our scheme, we can explain a little in the output. Similar to (Liznerski et al. 2020), our method outputs a vector with the same size as the input data where each dimension tells how much the noise is. Hence, if one dimension has a higher value, we can say the abnormality corresponds to this feature. We leave the detailed interpretability analysis in future research.

Extension to other data types Here, we provide insights for adapting noise evaluation to other data types like sequential, textual, and graph data. Adapting our approach to **images, time series, text, and graphs** is non-trivial but realizable, as each data type has distinct characteristics.

- **Image data** Images exhibit local smoothness or piecewise linear patterns across pixel values, so directly adding random noise (e.g., Gaussian) to pixel values is suboptimal. A better approach is to segment images into patches, sample contextual noise for each patch, and predict the average noise magnitude per segment.
- **Time series** Time series data, such as voice recordings, exhibit natural dependencies between adjacent samples, and anomalies may represent meaningful sequences. Pre-processing (e.g., ADbench(Han et al. 2022) for speech) can extract frequency features, transforming the series into a tabular format compatible with our model.
- **Textual data** Defining "noise" in textual data is non-trivial. A practical approach leverages recent large language models to extract latent features from text tokens. Using the CLS token or aggregation methods can yield

sentence- or document-level embeddings, allowing our noise evaluation model to operate in this latent space.

- **Graph data** Arbitrarily adding noise to the adjacency matrix is not feasible. One can add an extra virtual node(Cai et al. 2023) to capture the latent feature of the whole graph. Then, graph-level AD with our noise evaluation can be performed.

For all data types, instead of noise on raw data, we can apply diverse noise (e.g., Gaussian) to the latent features extracted by some (pre-trained) feature extractors (e.g. encoders) and then predict the noise amplitude. However, there is a risk that the feature extractor may fail to preserve the information of anomalies, which will lead to failure in detecting anomalies. A promising direction for further research involves training the feature extractor, a decoder, and the noise detector jointly, directly in the data space rather than the feature space.

Effectiveness on Dataset with Multiple Dense Regions

Figure 1 is just for illustration. Real-world normal data can come from several dense regions, which can be illustrated by a larger picture combining multiple patterns like Figure 1. Our method and theory are applicable to datasets spanning multiple clusters. For instance, $\tilde{\mathcal{D}}_H$ is the union of multiple hard-anomaly regions corresponding to different clusters.

We observe from real-world examples that our noise evaluation model remains effective when the dataset spans multiple clusters. For instance, in ADBench Letter and Vowels datasets, a t-SNE visualization in Figure 11 shows that Letter is composed of 5 clusters and Vowels has 3 dense regions. In Tables 5 and 6 of Appendix F, our method can obtain a high-performance score in terms of AUC (90.87%/99.03%) and F1 (67.79%/86.67%) on these datasets.

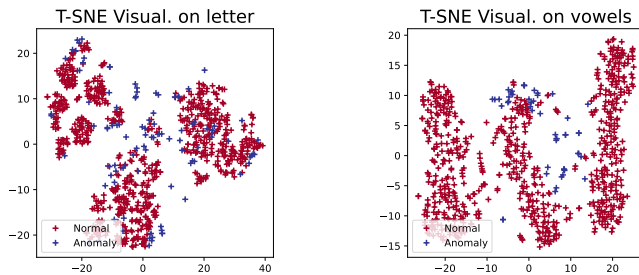


Figure 11: T-SNE Visual. of Letter and Vowels

Methods	DPAD	PLAD	NeuTraLAD	SCAD	AE	COPOD	DeepSVDD	ECOD	IForest	KNN	LOF	Ours-ResMLP	Ours-MLP
ALOI	64.48±0.4	53.16±1.4	53.54±0.9	53.97±0.5	54.96±0.1	51.51±0.1	51.44±1.3	53.07±0.1	53.78±0.4	66.72±0.2	74.69±0.2	60.59±0.7	59.95±0.7
anthyroid	90.91±2.0	95.09±2.8	81.18±2.6	86.52±2.2	84.66±1.5	77.59±0.2	56.97±2.6	79.11±0.2	91.39±1.4	93.6±0.4	94.78±0.3	97.02±0.5	96.78±0.3
backdoor	95.43±0.7	91.18±13.3	94.20±0.7	93.37±0.5	90.92±0.0	78.93±0.1	93.04±1.9	84.54±0.1	76.77±1.9	95.65±0.0	96.45±0.1	96.18±0.9	97.12±0.4
breastw	99.11±0.7	74.02±2.7	78.93±3.9	98.04±0.8	98.91±0.4	99.56±0.1	96.11±0.7	99.15±0.2	<u>99.47±0.1</u>	98.96±0.3	96.05±1.1	98.99±0.4	99.14±0.5
campaign	75.07±1.7	72.31±3.3	70.88±4.2	79.31±0.7	76.94±0.1	78.32±0.1	56.27±11.7	76.94±0.2	73.54±0.8	78.4±0.1	70.37±0.8	69.33±1.4	74.6±1.8
cardio	91.80±3.1	61.16±7.2	72.38±2.9	89.57±1.3	96.25±0.4	92.39±0.4	59.95±5.2	93.44±0.2	94.6±1.0	92.53±0.9	92.49±1.3	97.06±0.9	96.66±0.7
Cardio.	76.92±2.2	65.39±6.9	54.98±4.6	69.54±2.2	<u>83.35±1.0</u>	66.37±0.9	49.21±11.4	78.49±0.7	80.16±1.2	72.3±1.3	76.96±1.2	<u>85.04±3.2</u>	84.81±3.1
celeba	63.50±2.7	83.47±2.5	58.15±9.2	72.91±1.9	<u>79.8±0.1</u>	75.1±0.0	54.24±22.4	75.7±0.1	70.97±1.4	68.26±0.2	45.77±0.3	59.21±2.8	67.25±2.7
census	50.03±0.0	55.20±4.1	53.46±5.6	70.54±0.5	<u>70.74±0.1</u>	67.4±0.1	52.08±7.8	65.92±0.1	63.43±2.2	72.06±0.1	60.49±0.2	65.63±2.7	68.45±2.7
cover	60.42±10.8	87.66±0.5	71.61±10.4	96.68±1.2	95.53±0.9	88.41±0.0	36.57±5.8	92.01±0.1	85.07±3.1	98.6±0.1	<u>98.92±0.2</u>	99.69±0.1	99.64±0.1
donors	95.07±6.6	93.46±2.9	99.95±0.0	99.95±0.0	85.79±3.0	81.49±0.0	63.88±25.7	88.85±0.0	89.62±2.7	99.9±0.0	98.42±0.0	99.78±0.0	99.6±0.2
fault	73.06±1.0	61.51±4.8	69.39±1.3	<u>79.27±1.0</u>	55.26±1.2	45.13±0.8	50.6±6.2	46.79±0.6	65.84±1.6	79.83±1.1	65.91±1.7	72.27±1.3	73.49±1.2
fraud	89.73±5.6	96.29±0.7	91.59±3.5	94.71±0.4	95.38±0.0	94.74±0.0	85.15±6.7	94.97±0.0	94.88±0.4	96.17±0.1	77.85±3.4	94.44±0.3	95.11±0.3
glass	86.34±2.6	85.76±8.1	89.17±5.0	<u>92.23±2.5</u>	74.21±1.7	74.18±1.0	87.86±3.3	71.87±4.0	79.79±3.9	86.27±3.3	75.53±4.7	88.48±3.4	94.07±1.9
Hepatitis	77.67±8.9	73.08±3.9	51.70±11.2	69.7±6.2	79.63±4.6	80.89±2.4	68.44±6.0	74.22±2.6	79.02±5.8	<u>82.44±3.1</u>	80.93±3.2	81.73±3.0	87.75±2.3
http	<u>99.98±0.0</u>	99.99±0.2	79.82±32.1	99.86±0.3	99.94±0.0	99.15±0.0	21.43±43.1	97.86±0.0	99.14±0.6	99.96±0.0	93.8±0.4	99.99±0.0	99.99±0.0
InternetAds	89.99±0.9	87.44±1.5	88.94±1.0	88.76±1.4	87.84±0.3	67.92±0.7	89.91±0.4	67.63±0.5	44.39±2.1	88.41±0.7	88.79±0.6	92.49±0.7	92.36±0.7
Ionosphere	94.55±1.6	52.05±9.3	85.68±2.5	96.62±0.6	90.99±1.1	78.6±1.8	97.21±0.5	72.61±1.4	88.49±2.3	<u>97.42±0.7</u>	94.86±1.7	97.51±0.5	97.56±0.4
landsat	<u>74.67±2.2</u>	68.35±5.6	72.21±7.2	73.75±0.7	41.89±0.7	42.19±0.4	65.49±2.3	36.86±0.6	62.43±2.8	76.56±0.2	74.33±1.1	58.79±3.3	62.16±6.4
letter	72.34±2.2	70.19±5.5	<u>91.73±1.3</u>	94.36±0.8	52.36±0.4	55.53±0.9	62.58±2.8	57.18±0.7	62.87±2.4	88.12±0.9	86.85±1.2	90.87±0.8	90.71±1.2
Lympho.	98.71±0.2	71.43±38.5	44.48±9.6	99.37±0.9	98.36±2.4	99.73±0.4	97.95±1.5	<u>99.62±0.4</u>	99.2±0.9	98.41±2.5	99.19±0.7	98.95±1.9	98.64±1.9
magic.	77.94±4.6	78.46±2.7	71.09±2.0	80.05±0.3	70.18±0.9	68.24±0.3	63.49±0.6	63.55±0.2	77.16±1.1	83.31±0.3	83.29±0.2	88.05±0.4	88.59±0.4
mammo.	87.06±2.1	80.92±2.4	67.43±1.9	78.79±1.6	85.45±1.0	<u>90.61±0.1</u>	64.47±11.1	90.68±0.1	87.99±0.7	87.66±0.2	82.9±1.0	89.94±0.4	90.25±0.8
mnist	85.55±4.6	81.5±3.5	84.10±2.0	90.89±0.8	90.69±0.2	77.46±0.4	57.21±7.3	74.52±0.4	86.48±1.9	94.08±0.1	92.83±0.3	91.74±0.7	91.36±0.8
musk	99.99±0.0	91.51±7.8	100.00±0.0	100.00±0.0	100.00±0.0	94.56±0.2	100.00±0.0	95.74±0.3	95.6±3.0	100.00±0.0	100.00±0.0	100.00±0.0	100.00±0.0
optdigits	76.52±6.0	93.61±1.8	95.34±3.0	<u>96.49±0.9</u>	56.76±0.5	68.09±0.5	50.31±15.1	60.29±0.4	82.09±4.1	94.69±0.8	<u>97.33±0.4</u>	92.44±3.0	89.43±2.9
PageBlocks	96.12±0.4	85.77±3.0	90.79±0.6	94.52±0.7	94.75±0.2	87.45±0.3	75.82±9.9	91.32±0.2	92.35±0.5	96.3±0.2	<u>96.88±0.2</u>	97.78±0.9	97.76±0.4
pendigits	92.57±4.3	77.81±16.3	92.44±5.9	99.25±0.4	94.33±0.2	90.67±0.2	45.35±19.3	92.73±0.2	96.91±0.7	99.89±0.0	99.14±0.3	<u>99.64±0.1</u>	99.48±0.2
Pima	69.02±2.3	63.54±2.7	49.61±3.8	67.61±1.3	69.93±0.7	64.69±1.4	55.99±1.8	59.29±0.6	73.53±1.3	74.1±1.0	70.79±1.7	<u>71.2±1.7</u>	73.9±1.8
satellite	85.43±1.1	50.90±4.3	76.87±2.2	87.74±0.5	67.21±0.6	63.53±0.4	73.5±2.2	58.36±0.3	80.6±1.0	<u>87.59±0.2</u>	84.6±0.4	81.22±0.8	82.06±0.7
satimage	99.81±1.0	71.72±16.7	91.76±5.3	99.81±0.0	98.04±0.2	97.41±0.0	95.96±0.7	96.45±0.1	99.35±0.2	99.88±0.0	99.6±0.1	<u>99.87±0.0</u>	99.86±0.0
shuttle	99.96±0.0	99.89±0.1	99.92±0.1	99.99±0.0	99.43±0.1	99.45±0.0	99.66±0.0	99.3±0.0	99.58±0.1	99.89±0.0	99.98±0.0	99.99±0.0	99.98±0.0
skin	<u>99.48±0.2</u>	98.43±0.8	89.65±2.6	70.11±17.0	66.47±2.6	47.12±0.1	62.49±4.6	48.91±0.1	89.09±0.8	99.81±0.0	91.83±0.7	97.74±0.2	99.22±0.0
smtp	<u>96.00±2.0</u>	87.24±8.4	94.09±2.6	96.81±1.8	83.26±0.8	91.2±0.0	89.52±5.2	88.01±0.1	90.7±0.5	93.1±0.2	93.05±0.8	92.37±1.4	93.96±2.0
SpamBase	81.51±0.9	77.41±3.3	78.76±1.7	86.78±1.2	79.92±1.0	68.73±0.5	74.43±4.1	65.56±0.4	83.32±1.7	83.3±0.7	80.61±1.0	84.44±0.8	84.84±1.2
speech	<u>57.12±2.3</u>	58.68±2.2	55.43±5.1	56.55±2.0	47.29±0.4	49.18±0.3	50.83±1.4	47.29±0.3	48.18±1.5	48.62±0.9	49.64±0.6	55.98±1.9	58.13±2.2
Stamps	<u>95.06±1.2</u>	54.08±11.2	67.03±11.1	88.76±3.2	94.03±1.5	93.18±0.7	79.19±5.1	88.11±2.1	92.93±1.7	93.83±1.5	93.12±1.3	94.49±1.2	95.25±1.2
thyroid	97.95±1.5	97.39±2.9	74.34±7.5	94.95±1.6	98.71±0.1	93.93±0.2	80.64±10.1	97.73±0.2	98.97±0.2	<u>98.61±0.1</u>	96.31±1.1	97.74±0.5	98.49±0.4
vertebral	55.21±5.6	<u>58.72±13.8</u>	54.54±8.6	51.34±3.4	51.15±4.1	33.32±2.3	34.57±4.0	39.88±1.4	42.07±3.0	43.17±3.4	39.77±2.3	71.24±5.3	69.8±3.9
vowels	85.20±4.5	75.42±7.7	96.97±1.3	99.53±0.2	62.91±1.0	49.88±0.9	74.1±5.1	58.97±0.7	79.15±2.2	96.93±0.2	95.51±0.9	99.0±0.2	99.03±0.4
Waveform	68.80±5.4	83.53±1.2	72.88±1.2	68.54±3.3	65.22±0.5	73.42±0.7	60.59±5.5	60.04±0.3	73.0±2.0	<u>75.54±0.7</u>	75.34±0.8	75.29±1.6	72.38±3.0
WBC	98.66±0.3	48.87±8.2	63.30±8.5	95.61±0.3	99.04±0.6	99.58±0.3	93.32±1.9	99.53±0.3	<u>99.53±0.2</u>	98.98±0.5	97.3±1.5	98.41±0.8	99.17±0.7
WDBC	97.67±2.4	36.12±18.6	44.22±11.9	98.46±1.5	99.0±0.6	<u>99.3±0.3</u>	98.54±0.6	96.76±0.8	99.11±0.4	99.04±0.3	99.12±0.5	99.97±0.1	99.76±0.4
Wilt	62.79±3.2	99.03±3.1	65.44±2.4	74.16±4.6	37.9±2.0	34.43±0.4	45.2±2.0	39.31±0.4	48.08±2.2	64.24±1.8	69.54±3.0	<u>92.75±1.2</u>	94.3±0.9
wine	96.07±4.5	72.54±2.9	64.58±10.8	95.72±3.0	90.39±1.5	87.16±1.8	97.08±3.3	73.83±1.9	89.87±4.0	96.75±2.4	96.75±1.4	98.62±2.1	99.42±0.9
WPBC	56.61±2.0	62.95±6.8	53.51±3.5	46.67±3.3	49.65±3.2	50.32±2.8	49.9±2.3	48.21±1.4	50.83±3.1	52.99±2.8	52.54±2.8	59.79±1.8	61.25±2.4
yeast	51.35±3.1	<u>59.59±3.2</u>	52.10±3.4	50.95±2.2	42.18±0.8	38.28±0.5	45.06±3.1	44.15±1.5	40.66±1.0	45.06±1.3	46.92±1.4	58.54±1.8	57.53±2.0
mean rank	5.55	7.46	7.55	5.00	7.00	8.21	9.43	8.55	6.68	3.76	5.234	3.31	2.80
mean	82.75±15.6	75.40±18.3	74.47±17.9	84.44±18.0	78.46±18.7	74.60±19.6	68.37±22.7	74.15±19.3	79.82±17.5	85.91±15.6	83.58±16.5	87.07±14.3	87.89±13.69
p-value	0.0020	0.0000	0.0000	0.0325	0.0001	0.0000	0.0000	0.0000	0.0001	0.3128	0.0030	-	0.0088
p-value	0.0002	0.0000	0.0000	0.0039	0.0000	0.0000	0.0000	0.0000	0.0000	0.0716	0.0006	0.0088	-

Table 6: AUC (%) score of the proposed method compared with 8 baselines on 47 benchmark datasets. Each experiment is repeated 10 times with seed from 0 to 9, and mean value and standard deviation are reported. Mean is the average AUC score under all experiments. Mean rank is calculated out of 10 tested methods. Cardio. refers to Cardiocography. Mammo. refers to mammography. Lympho. refers to Lymphography.

Methods	DPAD	PLAD	NeuTraLAD	SCAD	AE	COPOD	DeepSVDD	ECOD	IForest	KNN	LOF	Ours-ResMLP	Ours-MLP
ALOI	14.15±0.2	9.40±0.8	7.79±0.6	11.79±0.6	8.02±0.1	5.26±0.2	6.92±0.5	5.42±0.1	6.72±0.6	15.87±0.3	20.16±0.8	14.32±0.6	13.53±0.7
annthyroid	63.97±3.6	76.22±6.8	42.06±5.0	54.68±0.9	51.02±1.7	31.6±0.3	21.99±3.0	39.03±0.3	56.69±2.3	63.73±1.6	67.87±1.6	75.59±2.6	74.74±1.6
backdoor	85.62±0.7	80.68±14.1	87.35±0.2	87.25±0.2	85.14±0.0	13.67±0.3	84.56±2.6	16.87±0.2	3.24±1.1	85.3±0.0	86.28±0.2	85.64±0.8	85.47±0.2
breastw	96.9±0.2	73.64±2.8	72.38±3.2	94.77±0.7	95.94±0.6	97.12±0.3	91.72±1.5	95.31±0.5	97.28±0.5	96.37±0.6	91.97±1.4	96.3±0.9	96.57±1.1
campaign	46.22±1.9	46.29±2.9	42.22±3.9	50.63±0.6	49.15±0.2	49.6±0.2	31.62±11.6	48.84±0.3	43.3±0.7	50.56±0.5	42.1±0.9	44.79±2.4	49.94±1.7
cardio	70.91±8.6	44.32±5.9	39.66±5.4	66.05±1.5	78.33±1.9	65.77±1.7	43.41±6.3	65.8±0.9	70.1±3.5	67.71±3.0	65.0±4.1	81.63±2.8	80.28±1.8
Cardiotocography	62.15±2.7	56.22±5.4	41.07±3.8	53.86±1.6	65.14±1.3	48.42±0.9	37.47±7.7	62.66±0.8	62.82±1.8	56.12±1.7	61.24±1.1	69.42±4.6	69.23±3.5
celeba	11.09±1.4	24.98±1.8	7.0±4.6	14.03±2.4	27.3±0.3	23.39±0.2	8.87±8.0	23.54±0.2	17.63±1.7	15.45±0.4	1.41±0.2	8.2±0.7	10.51±1.8
census	12.36±0.0	20.52±2.7	15.88±2.7	23.32±0.6	21.46±0.1	12.86±0.1	17.56±2.4	12.22±0.1	11.24±1.2	22.3±0.2	13.72±0.3	26.86±1.5	23.44±3.1
cover	25.37±24.1	33.63±2.3	19.27±16.1	65.25±9.7	20.16±5.0	18.17±0.2	2.72±2.9	23.54±0.2	11.33±1.9	74.89±0.3	79.03±2.1	91.38±0.3	89.97±1.0
donors	84.49±10.8	84.12±15.7	97.82±0.9	97.87±0.6	32.06±4.8	41.8±0.1	28.41±33.1	44.89±0.1	43.97±5.5	99.14±0.1	86.23±0.3	94.51±1.2	92.56±2.4
fault	67.7±0.7	60.22±3.2	67.19±1.1	72.62±0.9	54.87±0.6	47.77±0.7	50.88±4.9	48.65±0.1	61.68±1.7	73.13±1.0	64.19±0.6	68.03±0.9	68.3±1.0
fraud	51.91±10.9	78.66±6.0	45.05±21.9	65.13±2.4	32.44±0.6	44.72±0.5	52.48±14.7	38.86±0.6	30.77±4.5	47.12±2.6	0.0±0.0	60.19±9.2	78.68±1.6
glass	31.11±4.4	46.67±9.3	33.33±24.3	38.1±8.7	12.5±3.9	15.28±8.3	28.89±9.9	15.56±6.1	12.5±7.1	20.37±8.4	15.56±6.1	40.28±8.3	56.57±7.8
Hepatitis	60.0±12.3	51.92±3.8	29.23±10.2	36.26±8.6	51.92±6.8	53.85±0.0	44.62±10.0	40.0±3.4	50.96±10.0	55.13±5.8	52.31±3.4	65.38±4.1	69.23±4.9
http	99.68±0.1	99.22±0.8	1.9±0.1	84.57±31.3	92.09±0.8	1.99±0.0	9.91±17.7	2.2±0.0	13.9±15.3	99.56±0.1	0.05±0.0	99.25±0.5	99.57±0.2
InternetAds	82.34±0.7	79.46±1.8	81.63±1.2	81.17±2.1	79.93±1.3	50.75±0.5	81.3±1.5	50.43±0.3	24.12±2.6	80.12±2.1	80.65±0.7	84.24±1.4	84.19±0.8
Ionosphere	88.57±2.5	56.51±8.7	80.63±1.3	89.55±0.9	81.65±1.6	69.44±1.4	91.59±1.2	65.71±0.7	79.17±2.6	92.06±1.4	87.14±3.3	94.33±1.9	94.88±0.7
landsat	56.55±2.2	52.17±4.0	53.67±8.5	60.4±0.4	31.85±0.5	28.81±0.4	49.21±2.2	25.25±0.6	44.7±1.1	58.73±0.4	59.01±1.1	42.69±2.4	45.82±4.1
letter	35.2±5.3	44.50±6.6	58.6±6.5	63.17±3.9	13.88±0.8	14.0±1.4	26.0±2.5	16.0±1.2	16.88±3.4	46.83±3.9	48.4±3.2	67.79±2.0	65.45±2.8
Lymphography	83.33±0.0	70.83±23.3	0.0±0.0	88.89±13.6	79.17±24.8	93.75±8.6	70.0±18.3	90.0±9.1	83.33±8.9	80.56±24.5	83.33±16.7	86.9±17.9	75.76±21.6
magic	71.52±3.7	72.46±2.6	66.77±1.1	73.5±0.2	64.71±0.8	63.1±0.2	60.36±0.3	59.94±0.2	70.12±1.2	76.02±0.3	75.99±0.3	81.02±0.3	81.74±0.3
mammography	43.54±3.8	27.4±2.4	7.85±1.3	26.47±1.6	43.41±2.1	53.32±0.8	27.54±12.8	53.31±1.2	38.27±5.3	42.31±0.8	36.54±0.8	50.71±3.0	49.83±1.9
mnist	59.46±5.8	58.50±5.5	55.0±1.9	67.62±2.0	65.8±1.0	38.57±0.6	34.94±5.6	33.71±1.2	53.55±4.4	72.26±0.7	69.89±1.4	70.22±1.9	69.66±1.4
musk	98.76±2.0	75.26±22.1	100.0±0.0	100.0±0.0	100.0±0.0	48.97±1.1	100.0±0.0	55.05±1.6	54.12±16.7	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
optdigits	25.33±8.0	52.33±8.0	53.47±10.9	50.44±7.6	0.1±0.3	3.08±0.3	2.8±5.0	2.67±0.0	14.0±6.1	32.0±7.2	58.27±5.1	54.33±11.5	37.93±4.9
PageBlocks	80.63±1.6	68.24±4.4	62.67±1.8	75.42±2.0	68.96±0.5	48.97±0.9	60.31±9.5	57.22±0.7	63.46±1.8	78.99±0.9	81.29±0.5	83.37±1.0	83.94±1.4
pendigits	49.36±19.0	44.39±18.8	36.54±23.2	76.92±7.9	44.23±1.3	35.18±0.2	10.9±9.0	42.18±0.3	55.93±3.9	93.7±1.3	75.13±4.2	86.97±2.2	83.72±3.6
Pima	66.19±1.9	62.97±1.8	51.72±1.8	64.99±2.3	67.96±0.5	62.03±1.5	54.93±1.1	58.21±0.8	68.61±1.6	69.9±1.4	67.24±1.6	68.41±1.8	69.85±1.7
satellite	75.63±1.1	51.45±4.0	68.26±1.1	77.72±0.3	62.53±0.6	56.8±0.4	65.11±2.2	53.94±0.1	69.17±0.5	76.35±0.3	75.83±0.5	73.01±0.8	73.87±0.5
satimage	89.3±3.0	32.04±3.5	11.83±4.0	93.66±2.1	85.51±0.7	77.29±1.2	87.89±1.9	71.55±0.6	86.09±1.9	92.25±1.9	85.35±3.2	94.77±0.7	90.7±2.5
shuttle	98.17±0.2	98.63±0.3	98.15±1.0	99.01±0.1	95.95±0.2	95.96±0.1	97.56±0.5	91.55±0.3	94.72±1.8	98.03±0.2	98.26±0.2	98.56±0.2	98.48±0.2
skin	96.05±1.5	93.16±2.4	71.7±3.8	54.57±17.3	38.65±4.1	19.63±0.1	44.62±3.2	21.88±0.0	77.11±1.6	97.72±0.1	80.39±1.7	91.07±0.3	94.35±0.5
smtp	66.67±0.0	52.50±16.0	48.67±10.9	42.78±14.7	66.67±0.0	0.0±0.0	55.33±23.5	66.67±0.0	0.0±0.0	66.67±0.0	38.0±21.0	63.89±6.8	66.67±0.0
SpamBase	78.13±0.7	74.14±3.0	75.84±1.6	81.68±1.3	76.53±1.0	69.59±0.3	71.82±3.3	67.29±0.1	79.37±1.6	79.66±0.7	77.52±1.3	80.76±0.9	80.87±1.1
speech	8.85±2.5	10.25±4.5	7.54±4.2	5.25±2.7	4.92±0.0	3.28±0.0	3.93±1.9	4.92±0.0	5.33±2.1	4.92±1.0	6.23±0.7	12.02±2.0	13.11±1.5
Stamps	75.48±4.8	30.65±10.4	24.52±12.0	54.19±11.9	61.75±8.2	67.74±3.9	37.42±9.6	48.39±6.5	61.29±5.7	64.52±8.2	62.58±8.1	73.66±2.4	74.84±5.6
thyroid	76.77±5.5	74.19±6.5	7.96±7.0	61.72±7.2	76.5±1.1	30.24±0.4	47.96±17.6	61.94±1.4	81.32±2.9	73.84±1.6	58.06±5.0	70.25±5.9	76.24±5.2
vertebral	31.33±9.6	31.33±14.5	28.67±9.6	20.67±2.8	17.14±1.3	0.42±1.2	9.33±4.3	13.33±0.0	15.0±3.6	13.33±3.0	16.67±4.1	42.22±8.3	41.0±3.9
vowels	43.2±6.9	39.2±11.2	65.2±8.7	88.0±3.3	19.71±0.8	4.75±1.0	33.2±9.4	22.0±0.0	28.25±4.8	65.33±4.3	56.8±4.1	86.67±2.4	85.4±2.5
Waveform	20.6±2.9	57.00±3.4	47.8±1.6	19.75±4.4	9.14±0.7	9.5±1.1	15.8±2.5	7.4±0.5	10.75±1.8	29.5±2.7	28.8±1.1	27.0±3.3	20.2±2.5
WBC	82.0±4.0	30.00±33.9	12.0±7.5	57.5±5.0	85.71±5.3	87.5±7.1	54.0±11.4	88.0±4.5	86.25±5.2	83.33±8.2	68.0±11.0	81.67±7.5	83.0±9.5
WBBC	72.0±11.7	20.00±4.5	0.0±0.0	70.0±8.2	78.57±9.0	81.25±3.5	72.0±8.4	50.0±7.1	72.5±8.9	76.67±5.2	80.0±10.0	98.33±4.1	92.0±7.9
Wilt	14.86±5.4	89.88±3.4	28.87±5.4	31.23±5.5	2.95±0.7	1.61±0.1	3.97±1.4	5.68±0.2	1.9±0.4	2.01±0.5	17.98±4.0	62.52±5.7	68.6±2.9
wine	82.0±17.2	46.0±38.5	24.0±13.6	72.5±15.0	60.0±8.2	58.75±3.5	78.0±13.0	34.0±8.9	63.75±13.0	71.67±11.7	74.0±8.9	90.0±12.6	93.0±6.7
WPBC	48.09±3.2	48.94±5.3	41.7±5.1	32.98±2.1	36.17±5.2	31.38±4.5	37.02±2.4	36.17±0.0	35.64±2.2	35.46±2.6	38.3±3.4	50.35±2.6	51.06±2.0
yeast	52.54±2.7	58.54±2.3	52.24±3.2	52.07±1.9	44.24±0.8	41.72±0.6	47.61±1.2	46.43±1.2	43.44±1.0	47.4±1.1	49.03±1.1	57.66±1.3	57.22±1.6
mean rank	4.53	6.38	7.63	4.91	6.97	9.00	8.60	8.74	7.61	4.36	5.57	2.85	2.70
mean	60.34±27.4	55.1±25.1	44.1±29.1	60.64±26.2	51.95±29.0	40.82±28.1	44.56±29.5	42.22±25.1	46.22±29.1	62.66±27.8	57.06±28.4	69.09±24.7	69.4±24.5
p-value	0.0002	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0011	0.0000	-	0.68699
p-value	0.0001	0.0001	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0017	0.0002	0.68699	-

Table 7: F1 (%) score of the proposed method compared with 8 baselines on 47 benchmark datasets. Each experiment is repeated 10 times with seed from 0 to 9, and mean value and standard deviation are reported. Mean is the average F1 score under all experiments. Mean rank is calculated out of 10 tested methods. Cardio. refers to Cardiotocography. Mammo. refers to mammography. Lympho. refers to Lymphography.

Methods	DPAD	PLAD	NeuTraLAD	SCAD	AE	AnoGAN	COPOD	DeepSVDD	ECOD	IForest	KNN	LOF	Ours-ResMLP	Ours-MLP
abalone	97.23±2.7	97.12±2.8	96.96±2.9	97.05±2.8	97.17±2.6	97.13±2.7	97.01±2.7	96.99±2.9	97.00±2.7	97.22±2.6	97.09±2.8	97.15±2.7	97.38±2.5	97.39±2.5
arrhythmia	80.10±0.8	70.53±5.0	66.18±0.0	74.01±0.7	78.03±1.0	73.14±1.8	73.75±1.1	76.37±1.4	74.59±1.5	75.74±1.4	78.08±1.2	78.23±1.2	80.12±1.4	79.16±1.4
breastw	95.6±0.6	80.50±8.7	71.78±3.0	93.73±0.9	95.75±0.6	96.18±0.9	96.35±0.8	91.10±1.5	94.27±0.7	96.83±0.7	96.15±0.8	91.74±2.0	95.97±0.8	95.79±1.0
cardio	67.09±4.0	62.85±9.6	42.38±2.1	54.63±2.4	65.36±0.9	58.28±15.7	47.82±1.0	39.66±6.1	61.54±0.6	65.03±2.1	59.01±1.2	62.44±2.0	70.15±3.0	71.22±3.8
ecoli	96.09±1.7	89.38±10.5	86.32±9.4	94.05±2.2	95.77±2.1	95.10±2.1	88.91±4.2	93.42±2.5	87.55±6.2	95.71±1.7	95.87±1.8	95.73±1.8	96.47±1.5	96.61±1.5
glass	90.00±5.9	84.75±7.2	87.66±4.1	89.87±5.3	87.60±8.0	88.81±7.2	84.95±6.7	88.60±6.5	84.96±7.3	89.10±6.8	89.11±6.4	86.70±6.3	91.31±5.8	91.76±5.6
ionosphere	91.11±0.8	64.29±10.4	85.71±2.2	90.32±0.7	79.98±1.9	78.10±6.1	70.57±1.5	89.75±1.1	66.37±1.5	80.58±4.5	91.61±1.9	87.91±2.4	95.24±0.6	92.63±2.0
kdd	88.04±4.6	93.34±0.5	94.2±1.1	93.81±0.8	93.74±2.9	91.25±2.0	90.93±0.8	87.06±8.5	91.63±0.3	94.26±0.6	94.69±0.1	95.58±0.1	96.50±0.4	96.72±0.8
letter	98.86±0.3	98.01±0.1	98.99±0.2	99.60±0.1	98.67±0.3	98.47±0.3	98.05±0.1	98.77±0.2	98.05±0.1	98.98±0.2	99.4±0.2	98.86±0.3	99.55±0.1	99.57±0.1
lympho	77.05±2.7	68.85±1.6	58.36±3.4	74.59±4.4	75.91±3.0	71.80±4.8	61.61±2.2	73.77±4.6	62.76±2.7	73.36±3.8	75.78±3.3	75.55±3.2	80.56±3.3	81.73±2.9
mammo.	50.08±4.0	31.28±7.5	7.08±1.1	25.42±3.2	44.09±2.7	36.41±11.9	53.37±0.9	33.40±12.4	53.65±0.8	38.81±3.6	41.36±2.0	37.23±1.4	49.92±5.0	45.69±5.1
mulcross	99.95±0.1	99.67±0.6	25.41±26.5	99.95±0.1	100.0±0.0	98.52±2.0	66.00±0.1	100.0±0.0	74.58±0.2	99.11±0.6	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
musk	50.29±4.0	53.84±4.5	62.73±0.8	66.90±0.6	6.99±0.5	24.29±10.3	11.68±0.3	44.30±6.3	15.59±0.4	26.65±5.0	63.64±0.7	66.58±1.0	71.84±1.9	68.78±1.3
optdigits	96.66±1.7	91.73±2.5	92.83±2.2	98.64±0.9	97.95±1.1	95.84±1.7	92.18±1.4	96.95±1.4	92.04±1.4	98.13±1.1	98.60±0.8	98.63±0.8	98.80±0.8	98.91±0.8
pendigits	98.81±0.7	96.28±1.2	98.76±0.6	99.64±0.3	98.48±1.1	98.15±1.2	94.82±1.0	97.53±1.1	94.77±0.5	99.04±0.5	99.59±0.3	99.38±0.5	99.77±0.2	99.78±0.1
pima	67.31±1.3	64.55±3.3	51.34±3.2	65.07±1.1	67.72±1.2	67.01±3.1	62.69±1.2	56.96±2.6	58.12±0.9	68.94±1.4	70.06±1.0	67.24±1.2	69.03±2.0	69.78±2.1
satimage	94.76±2.3	83.49±7.0	83.19±7.0	94.76±1.9	95.05±2.4	94.75±2.6	89.43±2.9	86.66±5.2	87.21±3.5	96.05±2.0	95.65±2.6	93.89±3.0	95.57±2.1	95.50±2.2
seismic	34.12±3.0	35.29±4.1	17.18±1.9	29.06±2.9	28.12±1.4	29.56±1.7	30.07±1.5	24.89±6.2	29.9±0.3	29.36±1.5	30.22±2.0	16.12±1.8	31.18±1.3	34.03±1.7
shuttle	99.30±1.1	98.15±3.4	98.6±2.6	99.38±0.8	92.72±13.2	91.74±16.2	87.8±14.3	96.67±11.0	88.32±14.4	93.28±12.0	99.24±1.1	98.98±1.3	99.41±0.9	99.41±0.9
speech	7.54±2.5	9.84±3.8	4.92±4.1	5.41±2.2	4.64±0.6	3.28±2.3	3.28±0.0	4.31±2.3	4.92±0.0	4.37±2.3	5.76±0.9	6.84±1.0	13.35±1.8	11.94±1.6
thyroid	55.52±3.9	69.20±9.9	27.6±7.2	42.48±2.6	42.64±0.8	37.20±1.2	24.40±0.0	18.31±4.8	30.80±0.0	42.53±3.2	49.20±0.0	30.80±0.0	65.04±1.0	64.48±1.3
vertebral	90.19±1.1	81.27±6.0	81.14±0.8	86.33±1.5	89.97±0.9	89.62±2.3	87.24±0.8	85.17±1.5	82.63±0.9	89.38±1.1	88.93±0.7	89.51±1.1	92.04±0.7	91.50±0.7
vowels	40.4±8.5	48.67±34.5	62.91±8.0	86.20±2.6	19.94±0.8	9.60±7.4	4.44±0.8	33.54±8.8	21.20±1.0	28.33±4.1	64.86±4.2	57.49±3.5	87.71±3.5	86.57±1.9
wbc	89.25±1.7	66.67±12.1	57.26±8.7	86.89±1.4	88.97±0.9	87.45±3.6	79.53±1.0	83.61±4.0	62.74±1.1	89.9±1.1	88.54±1.0	88.96±0.9	92.65±1.3	92.05±1.6
wine	95.44±3.6	87.88±5.3	80.16±4.4	90.77±4.4	96.33±4.0	95.54±4.7	80.69±2.8	95.25±3.8	80.83±4.1	95.77±2.5	97.10±2.6	96.85±2.4	97.98±2.0	97.71±2.3
mean	92.67±15.1	89.00±17.4	87.95±20.5	92.8±15.7	91.37±18.3	88.58±21.4	87.83±19.4	90.39±18.5	87.87±18.3	91.85±17.3	93.43±14.8	92.6±16.1	94.23±14.0	94.18±13.8
mean rank	4.36	8.88	10.16	6.12	6.20	8.32	10.12	9.08	10.32	5.76	3.92	5.56	1.52	1.56
p-value	0.0292	0.0001	0.0000	0.0006	0.0190	0.0068	0.0008	0.0007	0.0002	0.0106	0.0023	0.0024	-	0.22
p-value	0.0363	0.0000	0.0000	0.0006	0.0198	0.0069	0.0008	0.0008	0.0002	0.0107	0.0030	0.0037	0.22	-

Table 8: F1 (%) score of the proposed method compared with 12 baselines on 25 benchmark datasets. Each experiment is repeated 10 times with seed from 0 to 9, and mean value and standard deviation are reported. Mean is the average F1 score under all experiments. Mean rank is calculated out of 10 tested methods. Mammo. refers to mammography.