# Anomaly Detection with Score Distribution Discrimination

Minqi Jiang
jiangmq95@163.com
AI Lab, Shanghai University of
Finance and Economics
Shanghai, China

Songqiao Han*
han.songqiao@shufe.edu.cn
AI Lab, Shanghai University of
Finance and Economics
Shanghai, China

Hailiang Huang*
hlhuang@shufe.edu.cn
AI Lab, Shanghai University of
Finance and Economics
Shanghai, China

## ABSTRACT

Recent studies give more attention to the anomaly detection (AD) methods that can leverage a handful of labeled anomalies along with abundant unlabeled data. These existing anomaly-informed AD methods rely on manually predefined score target(s), e.g., prior constant or margin hyperparameter(s), to realize discrimination in anomaly scores between normal and abnormal data. However, such methods would be vulnerable to the existence of anomaly contamination in the unlabeled data, and also lack adaptation to different data scenarios.

In this paper, we propose to optimize the anomaly scoring function from the view of score distribution, thus better retaining the diversity and more fine-grained information of input data, especially when the unlabeled data contains anomaly noises in more practical AD scenarios. We design a novel loss function called Overlap loss that minimizes the overlap area between the score distributions of normal and abnormal samples, which no longer depends on prior anomaly score targets and thus acquires adaptability to various datasets. Overlap loss consists of *Score Distribution Estimator* and *Overlap Area Calculation*, which are introduced to overcome challenges when estimating arbitrary score distributions, and to ensure the boundness of training loss. As a general loss component, Overlap loss can be effectively integrated into multiple network architectures for constructing AD models. Extensive experimental results indicate that Overlap loss based AD models significantly outperform their state-of-the-art counterparts, and achieve better performance on different types of anomalies.

## CCS CONCEPTS

• **Computing methodologies** → **Anomaly detection**; *Semi-supervised learning settings*.

## KEYWORDS

Anomaly Detection; Deep Learning; Neural Networks

---

*Corresponding author.

## 1 INTRODUCTION

Anomaly detection (AD) is the task of identifying unusual instances that deviate significantly from the majority of data, which has been applied in wide-ranging domains, such as social media analysis [41, 69], rare disease detection [51, 70], intrusion detection [30], and financial fraud detection [23, 50]. Previous research efforts [22, 35, 38, 55, 75] focus on unsupervised AD which does not require any labeled training data, but unsupervised methods lack any guidance of true anomalies [45, 48]. Therefore recent studies propose to learn valuable distinguishing features from a few labeled anomalies that may be identified by domain experts in practice, which is termed as the "anomaly-informed" AD methods [21, 48].

Current anomaly-informed methods, including semi- and weakly-supervised AD algorithms, mainly devise specific forms of loss functions to leverage such limited label information. These involve representation learning based minus loss in Unlearning [13], inverse loss in DeepSAD [56], and hinge loss in REPEN [45], as we summarized in Figure 1a~1c. In these methods, an anomaly score is generated based on the learned feature transformation of input data, such as the reconstruction error or the embedding distance. However, optimization in the representation space would lead to data-inefficient learning and suboptimal anomaly scoring [46, 48]. Therefore several works [47, 48, 73] fulfill an end-to-end learning fashion of anomaly score to obtain better performance, designing loss functions to map input instances to their corresponding anomaly score target(s), or to predefine a margin hyperparameter to realize the difference in anomaly scores between unlabeled samples and labeled anomalies, as shown in Figure 1d~1e. Nevertheless, such a predefined score target or margin would constrain the model's adaptability to different datasets, and further tuning these hyperparameters for realizing adaptation is often difficult, considering the scarcity of labeled data in practical AD scenarios.

In this paper, we aim to acquire an anomaly scoring function capable of realizing *adaptive anomaly score discrimination* for diverse data scenarios, thus addressing the above issues in previous anomaly-informed loss functions. Notably, we devise the Overlap loss that minimizes score distribution overlap between normal samples and anomalies, depending on the model itself to decide the suitable distribution of output anomaly scores. This kind of adaptability eliminates the dependency on predefined anomaly score targets. However, a non-trivial challenge is to estimate arbitrary distributions of anomaly scores, which are caused by the scarcity of labeled anomalies and anomaly noises in the unlabeled data. In the Overlap loss, we design a simple and effective method to estimate
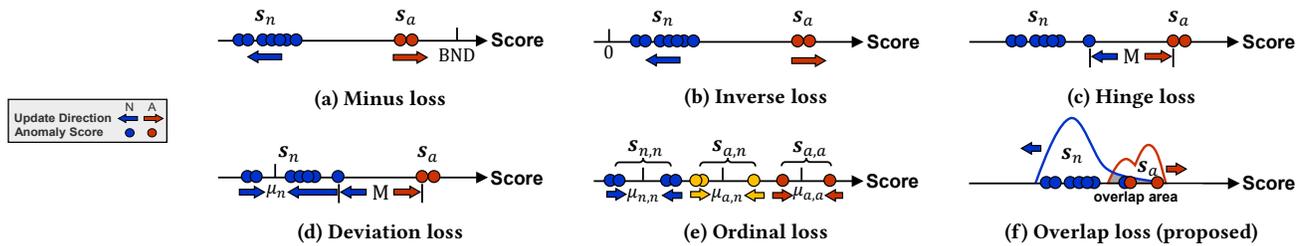
Figure 1: Loss function comparison. Unlike previous AD loss functions that suffer unbounded training loss or rely on specific anomaly score target(s) to guide model training, Overlap loss ensures a bounded training loss without requiring predefined target(s) by first estimating score distributions of normal and abnormal data, and further minimizing their overlap area.

the overlap area of arbitrary score distributions, while ensuring a correct order in the output anomaly scores and the boundness of training loss to better achieve stability in model training.

The main contributions of this paper can be summarized as follows: **(1)** We propose the Overlap loss for the AD community, which can achieve adaptive score distribution discrimination in input data, realizing sufficient global insight of anomaly scores in an end-to-end gradient update fashion. **(2)** We verify the effectiveness of the proposed Overlap loss on several network architectures covering both AD and classification tasks. Extensive results on 25 datasets suggest that the proposed Overlap loss could be served as a basis for further development in AD tasks. We open-source the proposed method, related codes, and all testing datasets for AD communities at https://github.com/Minqi824/Overlap. **(3)** We decouple the loss functions from several popular AD algorithms and analyze them in a unified framework, including embedding variation and network parameter changes. Moreover, we investigate the detection performance of different loss functions on various types of anomalies, therefore further exploring the pros and cons of these methods.

## 2 RELATED WORK

A desirable anomaly detection approach should produce not only a binary output (normal or abnormal) but also assign a degree of being an anomaly (anomaly score) to each observation [68]. Prior literature can be divided into two categories, i.e., AD algorithms without or with supervision. The former assumes that no labeled data is available during the model training stage and is proposed with different assumptions of data distribution [2], whereas the latter leverages a limited number of labeled samples which may be verified by some domain experts or automatic detecting systems.

**AD Algorithms without Supervision**. Typical anomaly detection methods are constructed for learning anomaly patterns in an unsupervised manner. These include shallow unsupervised models like CBLOF [22] and ECOD [35], or ensemble method Isolation Forest [38]. More recently, deep learning (DL) techniques like DeepSVDD [55] and GAN-based MO-GAAL [40] have been proposed for improving the performance of unsupervised AD tasks.

**AD Algorithms with Supervision**. Unsupervised methods can not achieve satisfactory performance in practical applications without the guidance of labeled data. Therefore several studies have also investigated utilizing partially labeled data to improve detection performance, which can be summarized into the following three categories:

(i) AD methods that are trained only on labeled normal samples, and detect anomalies that deviate from the normal representation learned in the training process [3, 4, 65, 66].

(ii) AD methods that additionally leverage a limited number of labeled anomalies. A common problem of the above methods using only normal samples is that many of the anomalies they identify are data noises or uninteresting data instances due to the lack of prior knowledge about the abnormal behaviors [45, 48]. This results in the development of semi- and weakly-supervised AD methods, which not only learn from numerous unlabeled data but also utilize limited information of labeled anomalies.

Among them, Unlearning [13] uses the minus loss form to provide an opposite direction of gradient update between the reconstruction error of the normal data and anomalies. DeepSAD [56] employs the inverse loss to penalize the inverse of the embedding distance such that the representation of anomalies must be mapped further away from the initial center of the hypersphere. REPEN [45] introduces a ranking model-based framework, which applies the hinge loss to encourage a distance separation of low-dimensional representation between normal samples and anomalies.

Several works are proposed to realize end-to-end learning of anomaly score, as they indicate that the above representation learning based AD methods would lead to data-inefficient learning and suboptimal detection performance [46]. Specified with the deviation loss, DevNet [48] leverages a prior probability and a margin hyperparameter to enforce significant deviations in anomaly scores between normal and abnormal data. FEAWAD [73] incorporates the DAGMM [75] network architecture with the deviation loss, for the better use of the information among hidden representation, reconstruction residual vector and reconstruction error transformed by the auto-encoder [7]. PReNet [47] formulates the scoring function as a pairwise relation learning task, where it defines three constant targets to enforce large margins among the anomaly scores of three types of instance pairs.

(iii) Fully-supervised methods are not specific for AD tasks in general [19]. Previous studies [6, 44] often use existing binary classifiers for this purpose such as Random Forest and MLP. One known risk of supervised methods is that ground truth labels maybe not necessarily accurate enough (i.e., there often exist some unlabeled anomaly noises in normal samples) to capture all types of anomalies, therefore these supervised methods may fail to detect unknown types of anomalies [21, 54].

In summary, some of the above anomaly-informed methods [13, 45, 56] perform an indirect representation learning of anomaly score, while other methods [47, 48, 73] mainly rely on predetermined training target(s) to realize the score discrepancy between the normal and abnormal data. Our proposed Overlap loss adaptively achieves the score discrimination from a distribution view, thus alleviating the need to define hyperparameter(s) as anomaly score target(s) in model training.

**Distribution Overlap**. Our idea is inspired by some recent studies of out-of-distribution (OOD) or multi-classification tasks in the CV field, whereas they usually consider the overlap of class distribution *only as a measurement* to describe the characteristics of datasets or to evaluate model quality [16, 26, 27, 39]. Considering research more closely related to our work, Magnet loss [53] is proposed to achieve local discrimination by penalizing class distribution overlap, as to realize explicit modeling of the distributions of different classes in representation space. Based on the entropy minimization principle [20], MA-DNN [11] minimizes the model entropy in the feature space and penalizes inconsistent network predictions at the class level. Nevertheless, these two methods are mainly devised for distance metric learning (DML) that presents optimization in the representation space. We propose to directly optimize distribution overlap in the anomaly scoring space to realize adaptive score distribution discrimination in input instances, which is tailored for the AD problem where there exist data noises and only a limited number of labeled anomalies.

## 3 METHODOLOGY

### 3.1 Problem Statement

Assume the training dataset $\mathcal{D} = \left\{ x_1^n, \ldots, x_k^n, \left( x_{k+1}^a, y_{k+1}^a \right), \ldots, \left( x_{k+m}^a, y_{k+m}^a \right) \right\}$ collects both unlabeled instances $\mathcal{D}_n = \left\{ x_i^n \right\}_{i=1}^k$ and a handful of labeled anomalies $\mathcal{D}_a = \left\{ \left( x_j^a, y_j^a \right) \right\}_{j=1}^m$, where $x \in \mathbb{R}^d$ represents the input feature and $y_j^a$ is the label of identified anomalies. Usually, we have $m \ll k$, since only limited prior knowledge of anomalies is available. Such data assumption is more practical for AD problems, and has been studied in recent works [24, 47, 48, 73]. Given such a dataset, our goal is to train a model, to effectively assign higher anomaly score for the abnormal data.

### 3.2 Overview of the Proposed Overlap Loss

Overlap loss first employs a *Score Distribution Estimator* for estimating the unknown probability density function (PDF) of the output anomaly scores in neural networks and then conducts *Overlap Area Calculation* between the anomaly score distributions of the unlabeled samples and labeled anomalies. Finally, Overlap loss minimizes the calculated overlap area of score distributions to provide the gradient for backpropagation in neural networks. The proposed Overlap loss fulfills the following properties: (i) the boundness of training loss for better convergence in the model training. (ii) eliminating explicit training target of anomaly score (e.g., constant or margin hyperparameter(s)) to enhance the model adaptability to different datasets. (iii) optimizing the entire anomaly score distribution, instead of pointwise optimization between the estimated anomaly scores and their corresponding targets.

## 3.3 Overlap Loss for Score Distribution Discrimination

In the following subsections, we illustrate two main parts of proposed Overlap loss: *Score Distribution Estimator* and *Overlap Area Calculation*, along with their corresponding basic ideas and challenges, as complements to our final solutions.

*3.3.1 Score Distribution Estimator.* Instead of pointwise optimization of the output anomaly scores, here we consider optimizing the anomaly score from a distribution view. Let $Q \in \mathbb{R}^M$ be the hidden representation space, an end-to-end anomaly scoring network $\phi(\cdot; \Theta) : x \mapsto \mathbb{R}$ can be defined as a combination of a feature representation learner $\psi(\cdot; \Theta_t) : x \mapsto Q$ and an anomaly scoring function $\eta(\cdot; \Theta_s) : Q \mapsto \mathbb{R}$, in which $\Theta = \{\Theta_t, \Theta_s\}$. If we denote the anomaly score of normal data as $\phi(x^n; \Theta) = s_n$ and that of abnormal data as $\phi(x^a; \Theta) = s_a$, a density estimator $f(\cdot)$ is then applied to estimate the PDFs of both $s_n$ and $s_a$ in a training batch.

A straightforward idea is to employ a prior distribution, e.g., the Gaussian distribution, as the score distribution estimator. Gaussian distribution inherits several good properties. For instance, the intersection point $c$ used for estimating the score distribution overlap in Figure 2a can be calculated by the following formula [25]:

$$c = \frac{\mu_a \sigma_n^2 - \sigma_a \left( \mu_n \sigma_a + \sigma_n \sqrt{(\mu_n - \mu_a)^2 + 2 \left( \sigma_n^2 - \sigma_a^2 \right) \log \left( \frac{\sigma_n}{\sigma_a} \right)} \right)}{\sigma_n^2 - \sigma_a^2}$$

(1)

where $\mu$ and $\sigma$ are their corresponding mean and variance of score distributions, respectively. The main challenge of this basic idea is that the number of labeled anomalies is usually too small to satisfy the Gaussian distribution assumption according to the central limit theorem [8], while enforcing the anomaly scores to follow this Gaussian prior would limit the representational ability of neural networks and further distort the anomaly scoring space, resulting in suboptimal performance.

To address the above challenges, we employ a score distribution estimator that is capable of estimating *arbitrary* distribution of output anomaly scores. In this paper, we use the non-parametric Kernel Density Estimation (KDE) method for estimating the arbitrary anomaly score distribution that may be caused by the scarcity of labeled data or the anomaly contamination in the unlabeled data. Actually, other differentiable density estimators can also be applied into our proposed Overlap loss.

If we denote the output anomaly score as $\phi(x; \Theta) = s$, the empirical cumulative distribution function (ECDF) can be defined as $\hat{F}_N(s) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{s_i \leq s}$, where $\mathbf{1}$ is the indicator function and $N$ is the number of partitions. $\hat{F}_N(s)$ is an unbiased estimator [12] of the cumulative distribution function (CDF) $F(s)$, and can be further used for estimating the PDF by the following equation:

$$\hat{f}(s) = \lim_{h \to 0} \frac{\hat{F}_N(s+h) - \hat{F}_N(s-h)}{2h} \approx \frac{1}{2Nh} \sum_{i=1}^N \left( \mathbf{1}_{s_i \leq s+h} - \mathbf{1}_{s_i \leq s-h} \right)$$

$$= \frac{1}{2Nh} \sum_{i=1}^N \left( \mathbf{1}_{s-h \leq s_i \leq s+h} \right) = \frac{1}{Nh} \sum_{i=1}^N \frac{1}{2} \mathbf{1} \left( \frac{|s - s_i|}{h} \leq 1 \right)$$

(2)

**(a) Gaussian score distribution overlap**　　**(b) Arbitrary score distribution overlap**　　**(c) Proposed score distribution overlap**
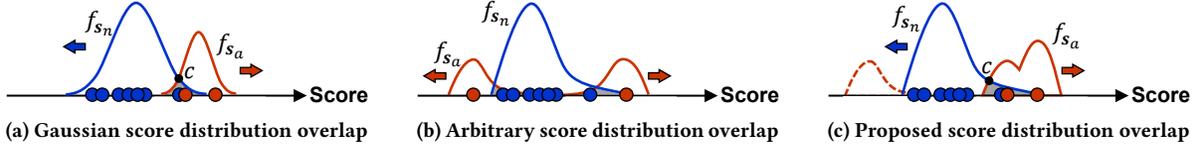
**Figure 2: Anomaly score distribution overlaps. (a) The prior assumption of Gaussian distributions limits the representational ability of neural networks. (b) Overlap of arbitrary score distributions leads to the disorder in anomaly scores. (c) The proposed Overlap loss minimizes the overlap area of arbitrary score distributions while ensuring correct order in anomaly scores.**

where $h$ is the bandwidth. If we denote the kernel function as $K(s) = \frac{1}{2}\mathbf{1}(s \leq 1)$, the estimated PDF can be rewritten as:

$$\hat{f}(s) = \frac{1}{Nh}\sum_{i=1}^{N} K\left(\frac{|s - s_i|}{h}\right) = \frac{1}{Nh}\sum_{i=1}^{N} K\left(\frac{s - s_i}{h}\right) \quad (3)$$

where $K(\cdot)$ is symmetric. We use Gaussian kernel in KDE, i.e., $K(s; h) \propto \exp\left(-\frac{s^2}{2h^2}\right)$, for estimating the unknown PDFs of anomaly scores, where the PDFs are further utilized to calculate the overlap area of score distributions, as described in the following subsection.

*3.3.2 Overlap Area Calculation.* Once we obtain the estimated score distributions (i.e., PDFs), the score distribution overlap can be calculated as the overlap area of PDFs between normal samples and the abnormal ones.

An optional method is to directly use the integral to approximate the score distribution overlap, as illustrated in Eq. 4. The overlap area between the PDFs of $s_n$ and $s_a$ is formulated as the integral of the one with the smaller probability density:

$$O(s_n, s_a) = \int_{\min(s_n, s_a)}^{\max(s_n, s_a)} \min\left(\hat{f}_{s_n}(t), \hat{f}_{s_a}(t)\right) dt \quad (4)$$

The main challenge of the above basic idea is that such a method *does not* necessarily guarantee a correct gradient update direction for anomaly scores, as illustrated in Figure 2b. The neural networks could minimize the overlap area of the score distributions, while mistakenly assigning lower anomaly scores for the anomalies (e.g., the left side in the score distribution of the anomalies in Figure 2b) instead of the normal ones. This problem can be remedied through the multi-task learning form by combining Eq.4 with a ranking loss term [61], as shown in Eq.5. However, although it ensures the order in anomaly scores, i.e., the anomaly scores of abnormal data should be further ranked higher than that of normal data, such a method may suffer from the difficult optimization problem in multi-task learning, sometimes leading to worse performance and data inefficiency compared to learning tasks individually [49, 64].

$$O(s_n, s_a) = \int_{\min(s_n, s_a)}^{\max(s_n, s_a)} \min\left(\hat{f}_{s_n}(t), \hat{f}_{s_a}(t)\right) dt + \max(0, s_n - s_a)$$

$$(5)$$

Our proposed Overlap loss aims to calculate the overlap area of *arbitrary* score distributions while ensuring the correct *order* in anomaly scores. We manage to acquire the intersection point $c$ of these arbitrary score distributions (see Figure 2c), and the score distribution overlap between $s_n$ and $s_a$ in a training batch can be further formulated as Eq.6, where $\hat{F}_{s_n}(\cdot)$ and $\hat{F}_{s_a}(\cdot)$ are the estimated CDF of normal and abnormal data, respectively.

$$O(s_n, s_a) = P(s_n > c) + P(s_a < c) = 1 - \hat{F}_{s_n}(c) + \hat{F}_{s_a}(c) \quad (6)$$

As shown in Figure 2c, the Overlap loss formulated in Eq.6 guarantees the order in output anomaly scores. A small overlap area with correct score order means a close to zero loss of $O(s_n, s_a)$. If the anomaly scores of abnormal data are smaller than that of normal data, $O(s_n, s_a)$ would penalize this disorder and be close to 2, since both $P(s_n > c)$ and $P(s_a < c)$ are close to 1, respectively. Moreover, $O(s_n, s_a)$ is naturally bounded to $[0, 2]$ due to the property of PDF.

However, for two arbitrary score distributions, we can not directly calculate the intersection point $c$ by the formula suitable for Gaussian distribution in Eq.1. Instead, we acquire the intersection point $c$ as the corresponding x value of the non-zero element of $d_k^s$ in Eq.7 for the arbitrary score distribution scenario, where $s_k$ is generated by the arithmetic sequence $s_k = \min(s_n, s_a) + (k - 1)\frac{\max(s_n, s_a) - \min(s_n, s_a)}{N}$ and $k = 1, \ldots, N$. In other words, we compare the PDF differences between two adjacent points of the score distributions $s_a$ and $s_n$, as shown in Figure 3.

$$d_k^s = \text{sgn}\left(\hat{f}_{s_a}(s_{k+1}) - \hat{f}_{s_n}(s_{k+1})\right) - \text{sgn}\left(\hat{f}_{s_a}(s_k) - \hat{f}_{s_n}(s_k)\right) \quad (7)$$

Figure 3 shows toy examples of calculating the intersection point(s) $c$. For most cases where there is only one intersection point between $\hat{f}_{s_n}(\cdot)$ and $\hat{f}_{s_a}(\cdot)$, $c$ is regarded as the x value of the sign change point of PDF differences, as shown in Figure 3a and 3d. Even if the two score distributions are far apart (see Figure 3b), we could still extend the x range of their PDFs and acquire $c$, as shown in Figure 3e. It is worth noting that for the case in Figure 3b, Overlap loss would reach its upper bound with an overlap area of 2 for penalizing the disorder in estimated anomaly scores, as illustrated in Eq.6. For the case where there exist multiple intersection points (shown in Figure3c and 3f), we randomly choose one of them as $c$. We show in the Appendix that the detection performance of this strategy is very close to that of ensembling different intersection points while improving the efficiency of model training.

After that, the integral of CDF $F_s(c)$ can be approximated via the trapezoidal rule [10] in Eq.8, where $\Delta s_k$ is adjusted based on the intersection point as $\Delta s_k = [c - \min(s_n, s_a)]/N$.

$$\hat{F}_s(c) = \int_{-\infty}^{c} \hat{f}_s(t)dt = \int_{\min(s_n, s_a)}^{c} \hat{f}_s(t)dt \approx \sum_{k=1}^{N} \frac{\hat{f}_s(s_k) + \hat{f}_s(s_{k+1})}{2}\Delta s_k$$

$$(8)$$

Based on the above notations, Overlap loss is defined as:

$$\mathcal{L}_{\text{Overlap}}(x \mid \Theta) = O(s_n, s_a) =$$

$$1 - \sum_{k=1}^{N} \frac{\hat{f}_{s_n}(s_{n,k}) + \hat{f}_{s_n}(s_{n,k+1})}{2}\Delta s_{n,k} + \sum_{k=1}^{N} \frac{\hat{f}_{s_a}(s_{a,k}) + \hat{f}_{s_a}(s_{a,k+1})}{2}\Delta s_{a,k}$$

$$(9)$$

**(a) PDFs w.r.t. the number of $c = 1$**

**(b) PDFs w.r.t. the number of $c = 0$**

**(c) PDFs w.r.t. the number of $c > 1$**

**(d) PDFs diff w.r.t. the number of $c = 1$**

**(e) PDFs diff w.r.t. the number of $c = 0$**

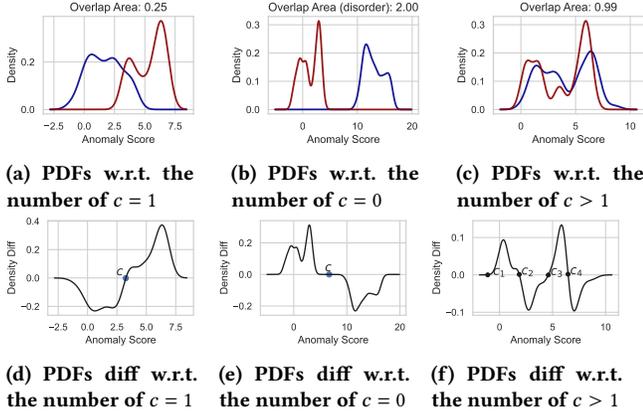**(f) PDFs diff w.r.t. the number of $c > 1$**

**Figure 3: Calculation of intersection point(s) for arbitrary anomaly score distributions of $f_{s_n}(\cdot)$ (blue) and $f_{s_a}(\cdot)$ (red). (a)~(c) correspond to the situations of one, zero, and multiple intersection point(s), respectively. (d)~(f) are their corresponding PDF differences.**

## 3.4 Network Architecture

Overlap loss is instantiated into an end-to-end neural network that consists of a feature representation layer $\psi(\cdot; \Theta_t)$ and a scoring layer $\eta(\cdot; \Theta_s)$. The BatchNorm layer is applied after the scoring layer to normalize the output anomaly scores. After that, score distributions of both normal data and anomalies are estimated by the KDE estimators, where their score distribution overlap is further calculated via the proposed Overlap loss, as shown in Figure 4.

We point out that the proposed Overlap loss can be effectively integrated into multiple popular network architectures, including the widely-used MLP and AutoEncoder in AD tasks, and some cutting-edge architectures like ResNet and Transformer in the classification tasks. Algorithm 1 provides detailed steps of instantiated models based on our proposed Overlap loss.

---

**Algorithm 1:** AD model instantiated by the Overlap loss

---

1 **Input**: Unlabeled instances $\mathcal{D}_n$, a limited number of identified anomalies $\mathcal{D}_a$

2 **Output**: Anomaly scores $s$

3 Initialize network parameters of both feature representation layer $\Theta_t$ and scoring layer $\Theta_s$

4 **for** *epoch=1:$n_{epoch}$* **do**

5    **for** *batch=1:$n_{batch}$* **do**

6      Randomly sample unlabeled instances $x_{batch}^n$ from $\mathcal{D}_n$ and labeled anomalies $x_{batch}^a$ from $\mathcal{D}_a$

7      (1) Acquire the anomaly scores with $\phi(x_{batch}^n; \Theta) = s_{batch}^n$ and $\phi(x_{batch}^a; \Theta) = s_{batch}^a$

8      (2) Use the KDE method to estimate the PDF $\hat{f}_{s_{batch}^n}(\cdot)$ and $\hat{f}_{s_{batch}^a}(\cdot)$ of the output anomaly scores

9      (3) Calculate the intersection point $c$ by Eq. 7

10      (4) Approximate the CDFs by the trapezoidal rule in Eq. 8

11      (5) Calculate and minimize the score distribution overlap $O\left(s_{batch}^n, s_{batch}^a\right)$ via Eq. 9

12      (6) Perform backpropagation and update network parameters $\Theta = \{\Theta_t, \Theta_s\}$

13    **end**

14 **end**

15 Output anomaly scores via learned scoring function $\phi(x; \Theta) = s$

---

# 4 EXPERIMENTS

## 4.1 Experiment Setting

**Datasets**. We apply 25 publicly available real-world datasets for model evaluation. These datasets include several domains such as disease diagnosis, speech recognition, and image identification. Detailed dataset description is illustrated in Appendix .1. For each dataset, 70% data is split as the training set and the remaining 30% as the testing set, where the same proportion of anomalies is kept by the stratified sampling. We discuss the model performance in Section 4.2.1 w.r.t. different ratios of labeled anomalies to all true anomalies $\gamma_l = m/(k + m)$ in the training set, where $m$ labeled anomalies are sampled from the entire anomaly data and the rest of $k$ instances remain unlabeled.

**Baselines**. We compare the proposed method with the following baselines[1], and summarize them according to their network architectures and levels of supervision, as is shown in Table 1.

- **Iforest** [38]. An ensemble of binary trees defines the anomaly score as the closeness of an individual instance to the root.
- **ECOD** [35]. A parameter-free method that estimates the empirical cumulative distribution of input features and regards tail probabilities as the anomaly score.
- **DeepSVDD** [55]. A neural network based model that describes the anomaly score as the distance of transformed embedding to the center of the hypersphere.
- **GANomaly** [3]. A GAN-based method that defines the reconstruction error of the input instance as the anomaly score.
- **DeepSAD** [56]. A deep semi-supervised one-class method that improves the unsupervised DeepSVDD.
- **REPEN** [45]. A neural network based model that leverages transformed low-dimensional representation for random distance-based detectors.
- **DevNet** [48]. A neural network based model that uses a prior probability to enforce the statistical deviation score of input instances.
- **PReNet** [47]. A neural network based model that defines a two-stream ordinal regression to learn the relation of instance pairs.
- **FEAWAD** [73]. A neural network based model that incorporates the network architecture of DAGMM [75] with the deviation loss of DevNet. We compare our proposed method with both its weakly- and fully-supervised versions.
- **ResNet** [18]. ResNet-like architecture turns out to be a strong baseline that is often missing in prior tabular AD tasks.
- **FTTransformer** [18]. A Transformer architecture implements with Feature Tokenizer. FTTransformer has been proven to be better than other DL solutions on tabular tasks.

**Metrics**. We evaluate the above models by two metrics: the AUC-ROC (Area Under Receiver Operating Characteristic Curve) and the AUC-PR (Area Under Precision-Recall Curve) values. We mainly report the AUC-PR results due to the space limit, and demonstrate the AUC-ROC results in the Appendix. We find that the results of these two metrics are generally consistent. Besides, we apply the

---

[1] Unlearning [13] is not included here since it is originally proposed for the time-series task, while we explore the Minus loss of Unlearning in Section 4.3. We do not include the results of DAGMM [75] for comparison as it may not converge on some datasets. We exclude the semi-supervised Dual-MGAN [34] method for comparison since it is too computationally expensive.
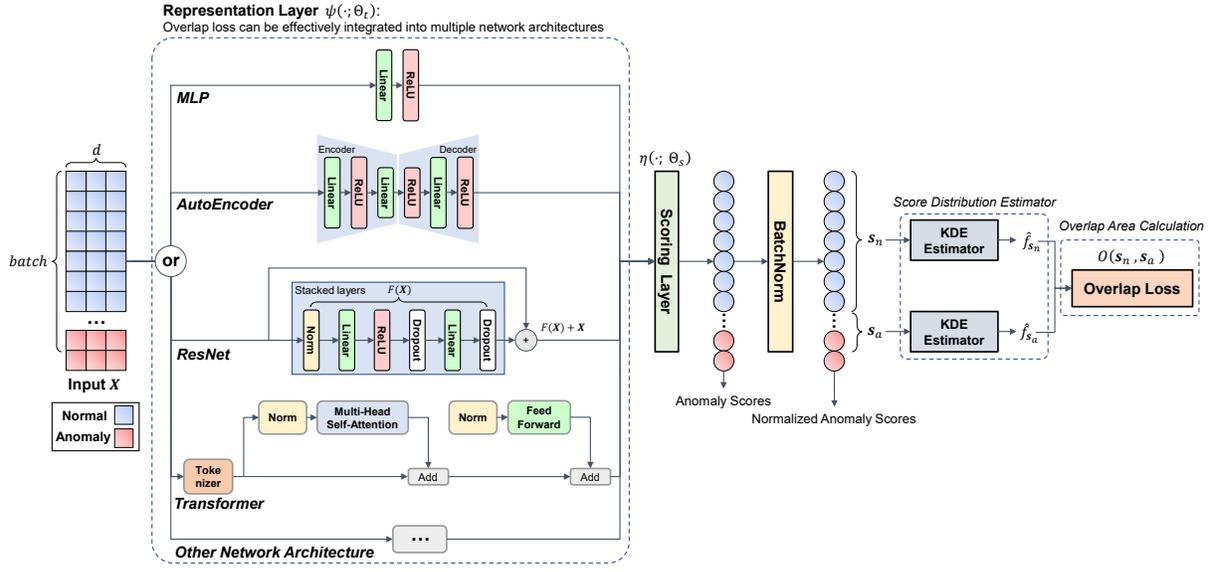
**Figure 4: AD model instantiated by the proposed Overlap loss, which consists of a representation layer $\psi(\cdot; \Theta_t)$ and a scoring layer $\eta(\cdot; \Theta_s)$ with batch normalization. The output anomaly scores are used for estimating the score distributions (PDFs) of normal samples $\hat{f}_{s_n}(\cdot)$ and that of anomalies $\hat{f}_{s_a}(\cdot)$ via the KDE estimators. Finally, the calculated overlap area of anomaly score distributions is minimized.**

pairwise Wilcoxon signed rank test [62] to examine the significance of proposed methods against its competitors.

**Training details** For the proposed Overlap loss based AD models, we use the SGD optimizer with 0.001 learning rate and 0.7 momentum. The weight decay is set to 0.01. The bandwidth $h$ in the KDE method is set to 1. The $N$ in Eq.7~9 is set to 1000 by default. We train the Overlap loss based MLP and AutoEncoder models (namely MLP-Overlap and AE-Overlap) for 20 epochs, where batch size of 256 is used. For ResNet and FTTransformer architectures, we train the Overlap loss based models (namely ResNet-Overlap and FTTransformer-Overlap) 100 training epochs just as their original paper. We provide the training details of compared baselines in Appendix .2, which are mainly according to their original papers. All the experiments are run on a Tesla V100 GPU accelerator.

## 4.2 Experimental Results

*4.2.1 Model Performance.* Table 1 shows the average model performance over 25 real-world datasets, and we report the full results in the Appendix of supplementary materials. Above all, we verify the effectiveness of the proposed Overlap loss on various network architectures, including MLP, AutoEncoder, ResNet, and Transformer. The Overlap loss based AD models generally outperform corresponding baselines w.r.t. the ratios of labeled anomalies $\gamma_l = 5\%$, $\gamma_l = 10\%$ and $\gamma_l = 20\%$.

Specifically, experimental results show that the MLP-Overlap achieves a relative improvement $\Delta$ Perf. of AUC-PR over its counterpart DevNet 2.89% and PReNet 1.82% w.r.t. $\gamma_l = 5\%$. These results indicate that compared to the current state-of-the-art weakly-supervised AD methods, Overlap loss is still more effective when only a handful of labeled anomalies (say 5% labeled anomalies)

are available in the training process, considering that such limited label information would bring challenges for estimating anomaly score distribution. Besides, we show that end-to-end AD methods, including DevNet, PReNet, and our MLP-Overlap, statistically outperform those unsupervised (e.g. Iforest) or semi-supervised representational learning (e.g., DeepSAD) AD methods, since end-to-end anomaly score learning can leverage the data much more efficiently than the two-step AD approaches [46, 48].

For $\gamma_l = 10\%$, the relative improvement of Overlap loss based model is more significant, as more labeled anomalies are beneficial for the Overlap loss to estimate more accurate anomaly score distributions, and thus to better measure the score distribution overlap. The MLP-Overlap $\Delta$ Perf. of AUC-PR over DevNet is 7.75% and 5.67% over PReNet w.r.t. $\gamma_l = 10\%$. Furthermore, we prove the superiority of the proposed Overlap loss on other network architectures such as AutoEncoder and ResNet. In terms of AUC-PR, the $\Delta$ Perf. of AE-Overlap over fully- and weakly-supervised FEAWAD is 28.04% and 9.29%, respectively, and $\Delta$ Perf. of ResNet-Overlap over ResNet is 56.30%, where all the relative improvements are significant at 1% significance level. Although FTTransformer has been proven to be a strong solution for tabular-based tasks [18], we still observe $\Delta$ Perf. of FTTransformer-Overlap over FTTransformer 5.50%~6.61% on AUC-PR.

*4.2.2 Runtime Analysis.* We show the model training time in Figure 5. This result shows that ECOD is the fastest algorithm as it treats each feature independently. Both MLP-Overlap and AE-Overlap are faster than their counterparts, since our methods need fewer training epochs while achieving better detection performance. For ResNet and FTTransformer (FTT) architectures, our methods are comparable to or relatively slower than the counterparts. This

**Table 1: Average AUC-PR performance over 25 real-world datasets. Each experiment is repeated 5 times. $\gamma_l$ stands for the ratio of labeled anomalies to all true anomalies in the training set. Δ Perf. shows the relative improvement of Overlap loss based models over their corresponding counterparts. \*\*\*, \*\* and \* denote statistical significance at 1%, 5% and 10% of Wilcoxon signed rank test, respectively. The best results are in bold.**

| Architecture | Model | Supervision | $\gamma_l = 5\%$ | | $\gamma_l = 10\%$ | | $\gamma_l = 20\%$ | |
|---|---|---|---|---|---|---|---|---|
| | | | AUC-PR | Δ Perf. | AUC-PR | Δ Perf. | AUC-PR | Δ Perf. |
| **Typical** | Iforest | Unsup | 0.389±0.295 | / | 0.389±0.295 | / | 0.389±0.295 | / |
| | ECOD | Unsup | 0.315±0.239 | / | 0.315±0.239 | / | 0.315±0.239 | / |
| | DeepSVDD | Unsup | 0.147±0.120 | / | 0.147±0.120 | / | 0.147±0.120 | / |
| | GANomaly | Semi | 0.297±0.191 | / | 0.296±0.195 | / | 0.306±0.201 | / |
| | DeepSAD | Semi | 0.506±0.253 | / | 0.601±0.275 | / | 0.675±0.284 | / |
| | REPEN | Weak | 0.560±0.300 | / | 0.603±0.308 | / | 0.639±0.306 | / |
| **MLP** | DevNet | Weak | 0.606±0.311 | +2.89% | 0.626±0.307 | +7.75%\*\* | 0.652±0.305 | +6.74%\* |
| | PReNet | Weak | 0.612±0.305 | +1.82% | 0.638±0.307 | +5.67%\* | 0.660±0.303 | +5.49% |
| | MLP-Overlap (ours) | Weak | **0.623±0.291** | / | **0.674±0.286** | / | **0.696±0.288** | / |
| **AutoEncoder** | FEAWAD | Sup | 0.509±0.269 | +28.04%\*\*\* | 0.620±0.270 | +12.05%\*\*\* | 0.678±0.270 | +5.17%\*\* |
| | FEAWAD | Weak | 0.596±0.286 | +9.29%\*\*\* | 0.645±0.293 | +7.71%\*\*\* | 0.682±0.283 | +4.56%\*\* |
| | AE-Overlap (ours) | Weak | **0.652±0.290** | / | **0.695±0.294** | / | **0.713±0.296** | / |
| **ResNet** | ResNet | Sup | 0.401±0.241 | +56.30%\*\*\* | 0.483±0.224 | +44.81%\*\*\* | 0.598±0.235 | +23.92%\*\*\* |
| | ResNet-Overlap (ours) | Weak | **0.627±0.297** | / | **0.699±0.289** | / | **0.742±0.283** | / |
| **Transformer** | FTTransformer | Sup | 0.594±0.299 | +5.50%\* | 0.644±0.308 | +6.61%\* | 0.691±0.305 | +5.65%\* |
| | FTTransformer-Overlap (ours) | Weak | **0.627±0.277** | / | **0.686±0.282** | / | **0.730±0.285** | / |

is mainly due to the fact that Overlap loss requires more training epochs for more complex network architectures (especially for FT-Transformer) than those simple architectures like MLP. Therefore we apply the same training strategy (100 epochs with early stopping) for ResNet and FTTransformer, as well as our Overlap loss based versions ResNet-Overlap and FTTransformer-Overlap. The extra training time is mainly caused by the calculation of Overlap loss, compared to the supervised binary cross entropy loss.



**Figure 5: Boxplot of model training time**

*4.2.3 Ablation Study.* In Table 2, we report the AUC-PR results of several basic methods mentioned in Section 3. We instantiate the basic method of Overlap-Gaussian by replacing the scoring layer $\eta\left(\cdot;\Theta_s\right)$ with the VAE [28] structure, where the anomaly scores of normal and abnormal data are sampled from their corresponding Gaussian distribution via the reparameterization trick [28]. The calculated intersection point $c$ of Eq.1 can be used for estimating score distribution overlap via Eq.6.

*First*, we observe that Overlap-Gaussian has the worst performance. This is because the scarceness of labeled anomalies makes its score distribution often present a certain arbitrariness, whereas the Gaussian assumption is detrimental to the representation of scoring function. *Second*, the disorder in anomaly scores leads to performance degradation in the Overlap-Arbitrary. Ranking loss term can be served as an effective way to guarantee the order in

anomaly scores, as the Overlap-Combined method significantly improves the AUC-PR performance. *Third*, the proposed Overlap loss outperforms all basic methods in most cases, since it can effectively estimate arbitrary score distributions of output anomaly scores while avoiding the score disorder problem that occurs in the Overlap-Arbitrary method. Compared to the Overlap-Combined method, Overlap-Proposed achieves better performance, probably because it realizes a unified loss function form, rather than a combination of two different loss parts. Besides, we observe that the multi-task loss form in the Overlap-Combined method fails in more complex network backbones like FTTransformer.

### 4.3 Further Exploration into AD Loss Functions

While most of the existing research focuses on proposing and evaluating specific models or architectural designs of AD methods [63], we manage to go a step further and directly compare different loss functions in the same network architecture. We introduce the decoupling methods in the Neural Architecture Search (NAS) problem [14, 37], where we mainly concern the design space of loss functions instead of other perspectives like architecture settings [33]. Such an analytical method could eliminate the effects of model configurations such as dropout and activation layers, while fully focusing on the role of loss functions (i.e., training objectives) in the anomaly detection tasks.

We decouple the loss functions in several popular AD models mentioned in Figure 1, including the Minus loss in Unlearning [13], Inverse loss in DeepSAD [56], Hinge loss in REPEN [45], Deviation loss in DevNet [48] and FEAWAD [73], Ordinal loss in PReNet [47], and our proposed Overlap loss, as shown in Table 3. For the consistency of comparison, we replace the original reconstruction error in Minus loss and the Euclidean distance of embedding in Inverse loss with the absolute anomaly score. A network with one-hidden-layer of 20 neurons is applied to ensure the comparability of different loss functions. The ReLU activation layer is employed

**Table 2: AUC-PR results of ablation studies. Overlap-Gaussian refers to the basic method mentioned in Section 3.3.1. Overlap-Arbitrary refers to the basic method of Eq.4. Overlap-Ranking isolates the ranking loss in Eq.5. Overlap-Combined corresponds to the combined loss form of both Overlap-Arbitrary and Overlap-Ranking as illustrated in Eq.5. Overlap-Proposed refers to the final solution in this paper.**

| Method | $\gamma_l = 5\%$ | | | | | $\gamma_l = 10\%$ | | | | | $\gamma_l = 20\%$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | VAE | MLP | AE | ResNet | FTT | VAE | MLP | AE | ResNet | FTT | VAE | MLP | AE | ResNet | FTT |
| Overlap-Gaussian | 0.159 | / | / | / | / | 0.158 | / | / | / | / | 0.158 | / | / | / | / |
| Overlap-Arbitrary | / | 0.349 | 0.376 | 0.341 | 0.356 | / | 0.376 | 0.416 | 0.407 | 0.339 | / | 0.397 | 0.410 | 0.426 | 0.352 |
| Overlap-Ranking | / | 0.535 | 0.558 | 0.539 | 0.623 | / | 0.605 | 0.620 | 0.625 | **0.687** | / | 0.657 | 0.657 | 0.691 | **0.731** |
| Overlap-Combined | / | **0.624** | 0.642 | 0.604 | 0.418 | / | **0.682** | **0.695** | 0.686 | 0.439 | / | **0.708** | **0.719** | 0.741 | 0.459 |
| Overlap-Proposed | / | 0.623 | **0.652** | **0.627** | **0.627** | / | 0.674 | **0.695** | **0.699** | 0.686 | / | 0.696 | 0.713 | **0.742** | 0.730 |

**Table 3: Summary of decoupled loss functions. The No Prior column indicates whether the prior anomaly score target is needed in the corresponding loss functions, e.g., the margin hyperparameter $M$ in the Hinge loss.**

| Loss | Formula | No Prior |
|---|---|---|
| Minus | $\mathcal{L} = |s_n| + \max(0, BND - |s_a|)$ | ✗ |
| Inverse | $\mathcal{L} = |s_n| + 1/|s_a|$ | ✓ |
| Hinge | $\mathcal{L} = \max(0, M + s_n - s_a)$ | ✗ |
| Deviation | $\mathcal{L} = |s_n| + \max(0, M - s_a)$ | ✗ |
| Ordinal | $\mathcal{L} = |s_{n,n} - \mu_{n,n}| + |s_{a,n} - \mu_{a,n}| + |s_{a,a} - \mu_{a,a}|$ | ✗ |
| Overlap | Eq.9 | ✓ |

in this network. We train the network for 200 epochs of 256 batch size and use the SGD optimizer with 0.01 learning rate and 0.7 momentum. The weight decay is set to 0.01. The hyperparameter $BND$ in the Minus loss is set to 5, and the hyperparameter $M$ in the Hinge and Deviation loss is set to 5. The anomaly scores in Deviation loss are normalized as Z-Score. Consistent with the original paper, we set $\mu_{n,n}$, $\mu_{n,n}$ and $\mu_{n,n}$ in the Ordinal loss to 0, 4, and 8, respectively. We first investigate different loss functions on 25 real-world datasets and then report their performances in detecting various types of anomalies.

*4.3.1 Exploration of AD Loss Functions on Real-world Datasets.* In this subsection, we analyze different loss functions on real-world datasets with respect to the following two perspectives: (i) Embedding transformation. The transformed embedding of input features [3, 4] can be seen as a visualization of representation layer variation for realizing the training objective. (ii) Network Parameter Changes. This is often discussed in the continual learning problem, where drastic changes in network parameters may suffer from the problem of catastrophic forgetting [5, 13, 29, 67]. Similarly, we investigate the network parameter changes of different loss function based AD models when achieving their corresponding training objectives.

**Embedding Transformation during Model Training**. We take the vowels dataset as an example to demonstrate the embedding transformation in the feature representation layer during the training process, as shown in Figure 7. Similar experimental results can be observed in other real-world datasets, and we provide these results in the Appendix of supplementary materials. Figure 7 indicates that the Deviation and Ordinal loss tend to seriously distort the embedding of original input data after a few training epochs. This is due to the fact that these two loss functions explicitly guide networks to map the anomaly score of each instance or pair to one

or more fixed score constants or a score margin, thus hindering the diversity of learned feature representation. Besides, the unlabeled normal samples are contaminated by the unlabeled anomalies, and defining an identical training target for these two types of data would limit the representational ability of the learned models.

Our proposed Overlap loss generates a relatively mild transformation of the input features. As mentioned before, Overlap loss based AD models can achieve superior detection performance, therefore we speculate that good detection performance *does not* always require an excessive transformation in the representation space, where the model can merely transform the embeddings that have the most impact on the score distribution discrimination and remain more fine-grained information of input data.

**Network Parameter Changes**. We show the results of network parameter changes on the 25 real-world datasets in Figure 6, where the sum of norms of parameter differences in each layer are calculated between the initialized model and its updated version as $\sum_l \|M - M_0\|_2^2$. The result indicates that compared to the other loss functions, the AD model based on our proposed Overlap loss inherits smaller network parameter changes. This result corresponds to the good properties of Overlap loss, where (i) Overlap loss is naturally bounded, avoiding drastic updating of anomaly scores (and network parameters). (ii) Overlap loss does not require the prior target of anomaly score, therefore reducing unnecessary scoring function updates in the training stage and being capable of adapting to different datasets with minimum adjustment of the score distribution.



**Figure 6: Network parameter changes in the training stage.**

*4.3.2 Exploration of AD Loss Functions on Different Types of Anomalies.* While extensive AD methods have been proven to be effective on real-world datasets, previous studies often neglect to discuss the pros and cons of AD methods regarding specific types of anomalies [17, 58]. In fact, public datasets often consist of a mixture of different types of anomalies. We follow [21, 58] to create realistic synthetic datasets based on the above 25 datasets by injecting four types (namely local, global, clustered, and dependency) of anomalies

(a) t-TSNE [60] plots of the input feature of vowels dataset.

(b) Minus

(c) Inverse

(d) Hinge

(e) Deviation

(f) Ordinal

(g) Overlap

Figure 7: Training loss along with the AUC-PR performance on testing set of different loss function based AD models, where the vowels dataset is specified for comparison. The transformed embeddings of the input feature are demonstrated, which corresponds to 5, 50, and 200 training epochs, respectively. See the additional results in Appendix.

to evaluate different loss functions. Appendix .5 provides detailed information of generated synthetic anomalies.

Table 4: Loss comparison on different types of anomalies.

| Loss | Local | Global | Clustered | Dependency |
|---|---|---|---|---|
| Minus | 0.255 | 0.822 | 0.992 | 0.369 |
| Inverse | 0.235 | 0.647 | 0.900 | 0.198 |
| Hinge | 0.271 | 0.853 | 0.996 | 0.413 |
| Deviation | 0.246 | 0.851 | 0.987 | 0.303 |
| Ordinal | 0.247 | 0.849 | 0.991 | 0.327 |
| Overlap | **0.439** | **0.929** | **0.998** | **0.571** |

Table 4 shows the AUC-PR results of loss function comparison on different types of anomalies. These results are consistent with the findings in [21], where the other loss functions devised for semi- or weakly-supervised AD algorithms perform relatively poorly on local and dependency anomalies. Unlike clustered anomalies, the partially labeled anomalies of local and dependency anomalies can not well capture all characteristics of specific types of anomalies, and learning such decision boundaries for separating normal and abnormal data is often challenging (see Figure A2g~A2k in the Appendix). Therefore, the incomplete label information may bias the learning process of these loss functions, which explains their relatively inferior performances compared to the Overlap loss.

In contrast, Overlap loss performs significantly better on local, global, and dependency anomalies, and achieves satisfactory results on clustered anomalies. For example, the average AUC-PR of Overlap loss on the local anomalies is 0.439, compared to the second-best Hinge loss of 0.271. Similar result can be observed for the dependency anomalies, where the AUC-PR of Overlap loss is 0.571, compared to the second-best Hinge loss of 0.413.

Such results verify that Overlap loss can effectively leverage the prior knowledge of both partial labels and anomaly types. That is to say, Overlap loss based AD models (e.g., ResNet-Overlap)

can achieve superior performance when only a limited number of labeled anomalies are available during the training stage. Furthermore, if one could get access to the valuable prior knowledge of anomaly types [21], Overlap loss can be served as an effective solution to learn the pattern of this specific type (e.g., dependency) of anomalies.

## 5 CONCLUSION

In this paper, we propose a novel loss function called Overlap loss for AD tasks. Overlap loss liberates the AD models from the predefined anomaly score targets, e.g., predefined constant or margin hyperparameter(s), thus adapting well to various datasets. By directly optimizing distribution overlap to realize score distribution discrimination, Overlap loss can retain more fine-grained information of input data, and also avoids dramatic changes in network parameters which may lead to overfitting or catastrophic forgetting problem. Extensive experimental results verify that the proposed Overlap loss can be effectively instantiated to different network architectures, including MLP, AutoEncoder, ResNet, and Transformer. Moreover, Overlap loss significantly outperforms other popular AD loss functions on various types of anomalies.

For the future, we plan to improve the optimization process of Overlap loss by leveraging more complex score distribution estimators, such as the Gaussian Mixture Model (GMM) [52]. Besides, we will extend our work to more general scenarios in weakly-supervised AD tasks [74], such as the inaccurate supervision [72] and inexact supervision [31, 59] problems.

## 6 ACKNOWLEDGMENTS

# REFERENCES

[1] Kjersti Aas, Claudia Czado, Arnoldo Frigessi, and Henrik Bakken. 2009. Pair-copula constructions of multiple dependence. *Insurance: Mathematics and economics* 44, 2 (2009), 182–198.

[2] Charu C Aggarwal. 2017. An introduction to outlier analysis. In *Outlier analysis*. Springer, 1–34.

[3] Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. 2018. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *Asian conference on computer vision*. Springer, 622–637.

[4] Samet Akçay, Amir Atapour-Abarghouei, and Toby P Breckon. 2019. Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection. In *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 1–8.

[5] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. 2018. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 139–154.

[6] Simon Duque Anton, Suneetha Kanoor, Daniel Fraunholz, and Hans Dieter Schotten. 2018. Evaluation of machine learning-based anomaly detection algorithms on an industrial modbus/tcp data set. In *Proceedings of the 13th international conference on availability, reliability and security*. 1–9.

[7] Pierre Baldi. 2012. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML workshop on unsupervised and transfer learning*. JMLR Workshop and Conference Proceedings, 37–49.

[8] Patrick Billingsley. 2008. *Probability and measure*. John Wiley & Sons.

[9] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. LOF: identifying density-based local outliers. In *SIGMOD*. 93–104.

[10] Richard L Burden, J Douglas Faires, and Annette M Burden. 2015. *Numerical analysis*. Cengage learning.

[11] Yanbei Chen, Xiatian Zhu, and Shaogang Gong. 2018. Semi-supervised deep learning with memory. In *Proceedings of the European conference on computer vision (ECCV)*. 268–283.

[12] Frederik Michel Dekking, Cornelis Kraaikamp, Hendrik Paul Lopuhaä, and Ludolf Erwin Meester. 2005. *A Modern Introduction to Probability and Statistics: Understanding why and how*. Vol. 488. Springer.

[13] Min Du, Zhi Chen, Chang Liu, Rajvardhan Oak, and Dawn Song. 2019. Lifelong anomaly detection through unlearning. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 1283–1297.

[14] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. 2019. Neural architecture search: A survey. *The Journal of Machine Learning Research* 20, 1 (2019), 1997–2017.

[15] Andrew Emmott, Shubhomoy Das, Thomas Dietterich, Alan Fern, and Weng-Keen Wong. 2015. A meta-analysis of the anomaly detection problem. *ArXiv* 1503.01158 (2015). https://arxiv.org/abs/1503.01158

[16] Eduardo Dadalto Camara Gomes, Florence Alberge, Pierre Duhamel, and Pablo Piantanida. 2021. Igeood: An Information Geometry Approach to Out-of-Distribution Detection. In *International Conference on Learning Representations*.

[17] Parikshit Gopalan, Vatsal Sharan, and Udi Wieder. 2019. Pidforest: anomaly detection via partial identification. *Advances in Neural Information Processing Systems* 32 (2019).

[18] Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. 2021. Revisiting deep learning models for tabular data. *Advances in Neural Information Processing Systems* 34 (2021), 18932–18943.

[19] Nico Görnitz, Marius Kloft, Konrad Rieck, and Ulf Brefeld. 2013. Toward supervised anomaly detection. *Journal of Artificial Intelligence Research* 46 (2013), 235–262.

[20] Yves Grandvalet and Yoshua Bengio. 2004. Semi-supervised learning by entropy minimization. *Advances in neural information processing systems* 17 (2004).

[21] Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. [n. d.]. ADBench: Anomaly Detection Benchmark. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

[22] Zengyou He, Xiaofei Xu, and Shengchun Deng. 2003. Discovering cluster-based local outliers. *Pattern recognition letters* 24, 9-10 (2003), 1641–1650.

[23] Waleed Hilal, S Andrew Gadsden, and John Yawney. 2021. A review of anomaly detection techniques and applications in financial fraud. *Expert Systems with Applications* (2021), 116429.

[24] Zongyuan Huang, Baohua Zhang, Guoqiang Hu, Longyuan Li, Yanyan Xu, and Yaohui Jin. 2021. Enhancing Unsupervised Anomaly Detection with Score-Guided Network. *arXiv preprint arXiv:2109.04684* (2021).

[25] Henry F Inman and Edwin L Bradley Jr. 1989. The overlapping coefficient as a measure of agreement between probability distributions and point estimation of the overlap of two normal densities. *Communications in Statistics-theory and Methods* 18, 10 (1989), 3851–3874.

[26] Dihong Jiang, Sun Sun, and Yaoliang Yu. 2021. Revisiting flow generative models for Out-of-distribution detection. In *International Conference on Learning Representations*.

[27] Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. 2021. Reversible Instance Normalization for Accurate Time-Series

Forecasting against Distribution Shift. In *International Conference on Learning Representations*.

[28] Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.).

[29] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526.

[30] Aleksandar Lazarevic, Levent Ertoz, Vipin Kumar, Aysel Ozgur, and Jaideep Srivastava. 2003. A comparative study of anomaly detection schemes in network intrusion detection. In *SDM*. SIAM, 25–36.

[31] Dongha Lee, Sehun Yu, Hyunjun Ju, and Hwanjo Yu. 2021. Weakly supervised temporal anomaly segmentation with dynamic time warping. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7355–7364.

[32] Meng-Chieh Lee, Shubhranshu Shekhar, Christos Faloutsos, T Noah Hutson, and Leon Iasemidis. 2021. Gen 2 Out: Detecting and Ranking Generalized Anomalies. In *Big Data*. IEEE, 801–811.

[33] Yuening Li, Zhengzhang Chen, Daochen Zha, Kaixiong Zhou, Haifeng Jin, Haifeng Chen, and Xia Hu. 2021. Autood: Neural architecture search for outlier detection. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2117–2122.

[34] Zhe Li, Chunhua Sun, Chunli Liu, Xiayu Chen, Meng Wang, and Yezheng Liu. 2022. Dual-MGAN: An Efficient Approach for Semi-supervised Outlier Detection with Few Identified Anomalies. *ACM Transactions on Knowledge Discovery from Data (TKDD)* (2022).

[35] Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and George Chen. 2022. Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions on Knowledge and Data Engineering* (2022).

[36] Boyang Liu, Pang-Ning Tan, and Jiayu Zhou. 2022. Unsupervised Anomaly Detection by Robust Density Estimation. In *AAAI*. AAAI Press, 4101–4108. https://ojs.aaai.org/index.php/AAAI/article/view/20328

[37] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. 2018. Progressive neural architecture search. In *Proceedings of the European conference on computer vision (ECCV)*. 19–34.

[38] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation forest. In *2008 eighth ieee international conference on data mining*. IEEE, 413–422.

[39] Siyan Liu, Pei Zhang, Dan Lu, and Guannan Zhang. 2021. PI3NN: Out-of-distribution-aware Prediction Intervals from Three Neural Networks. In *International Conference on Learning Representations*.

[40] Yezheng Liu, Zhe Li, Chong Zhou, Yuanchun Jiang, Jianshan Sun, Meng Wang, and Xiangnan He. 2019. Generative adversarial active learning for unsupervised outlier detection. *IEEE Transactions on Knowledge and Data Engineering* (2019).

[41] Jitendra Singh Malik, Guansong Pang, and Anton van den Hengel. 2022. Deep Learning for Hate Speech Detection: A Comparative Study. *arXiv preprint arXiv:2202.09517* (2022).

[42] Rafael Martinez-Guerra and Juan Luis Mata-Machuca. 2016. *Fault detection and diagnosis in nonlinear systems*.

[43] Glenn W Milligan. 1985. An algorithm for generating artificial test clusters. *Psychometrika* 50, 1 (1985), 123–127.

[44] Salima Omar, Asri Ngadi, and Hamid H Jebur. 2013. Machine learning techniques for anomaly detection: an overview. *International Journal of Computer Applications* 79, 2 (2013).

[45] Guansong Pang, Longbing Cao, Ling Chen, and Huan Liu. 2018. Learning representations of ultrahigh-dimensional data for random distance-based outlier detection. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 2041–2050.

[46] Guansong Pang, Choubo Ding, Chunhua Shen, and Anton van den Hengel. 2021. Explainable Deep Few-shot Anomaly Detection with Deviation Networks. *arXiv preprint arXiv:2108.00462* (2021).

[47] Guansong Pang, Chunhua Shen, Huidong Jin, and Anton van den Hengel. 2019. Deep weakly-supervised anomaly detection. *arXiv preprint arXiv:1910.13601* (2019).

[48] Guansong Pang, Chunhua Shen, and Anton van den Hengel. 2019. Deep anomaly detection with deviation networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 353–362.

[49] Emilio Parisotto, Lei Jimmy Ba, and Ruslan Salakhutdinov. 2016. Actor-Mimic: Deep Multitask and Transfer Reinforcement Learning. In *ICLR (Poster)*.

[50] Tahereh Pourhabibi, Kok-Leong Ong, Booi H Kam, and Yee Ling Boo. 2020. Fraud detection: A systematic literature review of graph-based anomaly detection approaches. *Decision Support Systems* 133 (2020), 113303.

[51] Theodoros Rekatsinas, Saurav Ghosh, Sumiko R Mekaru, Elaine O Nsoesie, John S Brownstein, Lise Getoor, and Naren Ramakrishnan. 2015. Sourceseer: Forecasting rare disease outbreaks using multiple data sources. In *Proceedings of the 2015 SIAM International Conference on Data Mining*. SIAM, 379–387.

[52] Douglas A Reynolds. 2009. Gaussian mixture models. *Encyclopedia of biometrics* 741, 659-663 (2009).

[53] Oren Rippel, Manohar Paluri, Piotr Dollar, and Lubomir Bourdev. 2015. Metric learning with adaptive density discrimination. *arXiv preprint arXiv:1511.05939* (2015).

[54] Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. 2021. A unifying review of deep and shallow anomaly detection. *Proc. IEEE* 109, 5 (2021), 756–795.

[55] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep one-class classification. In *International conference on machine learning*. 4393–4402.

[56] Lukas Ruff, Robert A. Vandermeulen, Nico Görnitz, Alexander Binder, Emmanuel Müller, Klaus-Robert Müller, and Marius Kloft. 2020. Deep Semi-Supervised Anomaly Detection. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

[57] Jonas Soenen, Elia Van Wolputte, Lorenzo Perini, Vincent Vercruyssen, Wannes Meert, Jesse Davis, and Hendrik Blockeel. 2021. The Effect of Hyperparameter Tuning on the Comparative Evaluation of Unsupervised Anomaly Detection Methods. In *Proceedings of the KDD'21 Workshop on Outlier Detection and Description*. Outlier Detection and Description Organising Committee, 1–9.

[58] Georg Steinbuss and Klemens Böhm. 2021. Benchmarking unsupervised outlier detection with realistic synthetic data. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 15, 4 (2021), 1–20.

[59] Waqas Sultani, Chen Chen, and Mubarak Shah. 2018. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6479–6488.

[60] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).

[61] Xinshao Wang, Yang Hua, Elyor Kodirov, Guosheng Hu, Romain Garnier, and Neil M Robertson. 2019. Ranked list loss for deep metric learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5207–5216.

[62] Robert F Woolson. 2007. Wilcoxon signed-rank test. *Wiley encyclopedia of clinical trials* (2007), 1–3.

[63] Jiaxuan You, Zhitao Ying, and Jure Leskovec. 2020. Design space for graph neural networks. *Advances in Neural Information Processing Systems* 33 (2020), 17009–17021.

[64] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. 2020. Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems* 33 (2020), 5824–5836.

[65] Houssam Zenati, Chuan Sheng Foo, Bruno Lecouat, Gaurav Manek, and Vijay Ramaseshan Chandrasekhar. 2018. Efficient gan-based anomaly detection. *arXiv preprint arXiv:1802.06222* (2018).

[66] Houssam Zenati, Manon Romain, Chuan-Sheng Foo, Bruno Lecouat, and Vijay Chandrasekhar. 2018. Adversarially learned anomaly detection. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 727–736.

[67] Friedemann Zenke, Ben Poole, and Surya Ganguli. 2017. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*. PMLR, 3987–3995.

[68] Liangwei Zhang, Jing Lin, and Ramin Karim. 2018. Adaptive kernel density-based anomaly detection for nonlinear systems. *Knowledge-Based Systems* 139 (2018), 50–63.

[69] Jun Zhao, Xudong Liu, Qiben Yan, Bo Li, Minglai Shao, and Hao Peng. 2020. Multi-attributed heterogeneous graph convolutional network for bot detection. *Information Sciences* 537 (2020), 380–393.

[70] Yue Zhao, Xiyang Hu, Cheng Cheng, Cong Wang, Changlin Wan, Wen Wang, Jianing Yang, Haoping Bai, Zheng Li, Cao Xiao, et al. 2021. SUOD: Accelerating large-scale unsupervised heterogeneous outlier detection. *MLSys* 3 (2021), 463–478.

[71] Yue Zhao, Zain Nasrullah, and Zheng Li. 2019. PyOD: A Python Toolbox for Scalable Outlier Detection. *Journal of Machine Learning Research* 20 (2019), 1–7.

[72] Yue Zhao, Guoqing Zheng, Subhabrata Mukherjee, Robert McCann, and Ahmed Awadallah. 2022. Admoe: Anomaly detection with mixture-of-experts from noisy labels. *arXiv preprint arXiv:2208.11290* (2022).

[73] Yingjie Zhou, Xucheng Song, Yanru Zhang, Fanxing Liu, Ce Zhu, and Lingqiao Liu. 2021. Feature encoding with autoencoders for weakly supervised anomaly detection. *IEEE Transactions on Neural Networks and Learning Systems* (2021).

[74] Zhi-Hua Zhou. 2018. A brief introduction to weakly supervised learning. *National science review* 5, 1 (2018), 44–53.

[75] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng Chen. 2018. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International conference on learning representations*.

# APPENDIX

## .1 Details of Dataset Description

We show the detailed description of 25 publicly available real-world datasets in Table A1, which include several domains such as disease diagnosis, speech recognition, and image identification.

**Table A1: Dataset description.**

| Dataset | $N$ | $D$ | #anomalies | #anomaly ratio (%) |
|---|---|---|---|---|
| ALOI | 49534 | 27 | 1508 | 3.04 |
| annthyroid | 7200 | 6 | 534 | 7.42 |
| Cardiotocography | 2114 | 21 | 466 | 22.04 |
| fault | 1941 | 27 | 673 | 34.67 |
| http | 567498 | 3 | 2211 | 0.39 |
| landsat | 6435 | 36 | 1333 | 20.71 |
| letter | 1600 | 32 | 100 | 6.25 |
| magic.gamma | 19020 | 10 | 6688 | 35.16 |
| mammography | 11183 | 6 | 260 | 2.32 |
| mnist | 7603 | 100 | 700 | 9.21 |
| musk | 3062 | 166 | 97 | 3.17 |
| optdigits | 5216 | 64 | 150 | 2.88 |
| PageBlocks | 5393 | 10 | 510 | 9.46 |
| pendigits | 6870 | 16 | 156 | 2.27 |
| satellite | 6435 | 36 | 2036 | 31.64 |
| satimage-2 | 5803 | 36 | 71 | 1.22 |
| shuttle | 49097 | 9 | 3511 | 7.15 |
| skin | 245057 | 3 | 50859 | 20.75 |
| SpamBase | 4207 | 57 | 1679 | 39.91 |
| speech | 3686 | 400 | 61 | 1.65 |
| thyroid | 3772 | 6 | 93 | 2.47 |
| vowels | 1456 | 12 | 50 | 3.43 |
| Waveform | 3443 | 21 | 100 | 2.90 |
| Wilt | 4819 | 5 | 257 | 5.33 |
| yeast | 1484 | 8 | 507 | 34.16 |

## .2 Additional Training Details

For unsupervised baselines, Iforest, ECOD, and DeepSVDD are built using the PyOD [71] library. Labeled anomalies are combined with unlabeled data for constructing the validation set, in order to tune the hyperparameters of these unsupervised methods via the grid search method, since tuning their hyperparameters on a small validation set often yields better performance than using the default settings [57]. Table A2 shows the hyperparameter grids, where ECOD is not considered since it is a parameter-free method.

**Table A2: Hyperparameter grid of the unsupervised models.**

| Model | Hyperparameter |
|---|---|
| Iforest | n_estimators: [10, 50, 100, 500] |
| DeepSVDD | epochs: [20, 50, 100, 200] |

We replace the convolutional layer in the original GANomaly with the dense layer for evaluating it on the tabular data, where the hidden size of the encoder-decoder-encoder structure of GANomaly

is set to half of the input dimension. We realize the PReNet in Py-Torch as we do not find the open-source codes, and set the hyperparameters in PReNet according to its original paper. Other models are built based on their corresponding source codes. If not specified, we train these baseline models according to their default hyperparameters mentioned in the original papers.

## .3 Experimental Results of Selecting Intersection Points

We demonstrate the experimental results of different intersection point selection strategies in Table A3. The detection performance of a randomly selected intersection point is very close to that of ensembling all the calculated intersection points w.r.t. different ratios of labeled anomalies $\gamma_l = 5\%$, $\gamma_l = 10\%$ and $\gamma_l = 20\%$. This is due to the fact that the calculation of the overlap area repeats many times (epochs×batchsize), which is essentially similar to the average results of the ensemble strategy. Moreover, random sampling of the intersection points can improve computational efficiency since the overlap area only needs to be estimated once in a training batch.

**Table A3: AUC-ROC and AUC-PR results of different intersection point selection strategies. MLP-Overlap corresponds to the default strategy that randomly chooses one of the intersection points for estimating the overlap area via Eq.9. MLP-Overlap-E refers to the ensemble strategy by taking the average of the overlap areas calculated based on each intersection point.**

(a) AUC-ROC results.

| | $\gamma_l = 5\%$ | $\gamma_l = 10\%$ | $\gamma_l = 20\%$ |
|---|---|---|---|
| MLP-Overlap | 0.847±0.145 | 0.880±0.132 | 0.893±0.133 |
| MLP-Overlap-E | 0.849±0.146 | 0.879±0.132 | 0.894±0.131 |

(b) AUC-PR results.

| | $\gamma_l = 5\%$ | $\gamma_l = 10\%$ | $\gamma_l = 20\%$ |
|---|---|---|---|
| MLP-Overlap | 0.623±0.291 | 0.674±0.286 | 0.696±0.288 |
| MLP-Overlap-E | 0.626±0.291 | 0.673±0.285 | 0.694±0.287 |

## .4 Detailed Experimental Results

We show the AUC-ROC results of model performance in Table A4 and ablation study in Table A5, corresponding to Section 4.2.1 and 4.2.3 in the main paper, respectively. These experimental results are basically consistent with the main paper. Besides, we show the detailed results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 5\%$, $\gamma_l = 10\%$ and $\gamma_l = 20\%$ in Table A6~A11. The best performing method(s) is marked in **bold**.

**Table A4: Average AUC-ROC performance over 25 real-world datasets. Each experiment is repeated 5 times. $\gamma_l$ stands for the ratio of labeled anomalies to all true anomalies in the training set. $\Delta$ Perf. shows the relative improvement of Overlap loss based models over their corresponding counterparts. ***, ** and * denote statistical significance at $1\%$, $5\%$ and $10\%$ of Wilcoxon signed rank test, respectively. The best results are in bold.**

| Architecture | Model | Supervision | $\gamma_l = 5\%$ | | $\gamma_l = 10\%$ | | $\gamma_l = 20\%$ | |
|---|---|---|---|---|---|---|---|---|
| | | | AUC-ROC | $\Delta$ Perf. | AUC-ROC | $\Delta$ Perf. | AUC-ROC | $\Delta$ Perf. |
| **Typical** | Iforest | Unsup | 0.737±0.187 | / | 0.737±0.187 | / | 0.737±0.187 | / |
| | ECOD | Unsup | 0.701±0.208 | / | 0.701±0.208 | / | 0.701±0.208 | / |
| | DeepSVDD | Unsup | 0.504±0.028 | / | 0.504±0.028 | / | 0.504±0.028 | / |
| | GANomaly | Semi | 0.655±0.162 | / | 0.648±0.153 | / | 0.665±0.152 | / |
| | DeepSAD | Semi | 0.823±0.142 | / | 0.859±0.136 | / | 0.888±0.129 | / |
| | REPEN | Weak | 0.810±0.166 | / | 0.832±0.165 | / | 0.848±0.163 | / |
| **MLP** | DevNet | Weak | 0.842±0.148 | +0.57% | 0.861±0.135 | +2.16% | 0.873±0.129 | +2.40% |
| | PReNet | Weak | 0.846±0.146 | +0.18% | 0.866±0.132 | +1.61% | 0.876±0.127 | +1.97% |
| | MLP-Overlap (ours) | Weak | **0.847±0.145** | / | **0.880±0.132** | / | **0.893±0.133** | / |
| **AutoEncoder** | FEAWAD | Sup | 0.771±0.211 | +11.72%*** | 0.849±0.133 | +4.34%*** | 0.876±0.133 | +2.42%*** |
| | FEAWAD | Weak | 0.808±0.154 | +6.60%*** | 0.848±0.145 | +4.51%*** | 0.876±0.129 | +2.43%*** |
| | AE-Overlap (ours) | Weak | **0.862±0.144** | / | **0.886±0.137** | / | **0.897±0.132** | / |
| **ResNet** | ResNet | Sup | 0.651±0.158 | +28.48%*** | 0.736±0.124 | +19.74%*** | 0.816±0.127 | +11.03%*** |
| | ResNet-Overlap (ours) | Weak | **0.836±0.146** | / | **0.882±0.134** | / | **0.906±0.122** | / |
| **Transformer** | FTTransformer | Sup | 0.827±0.159 | +3.00%** | 0.859±0.146 | +1.80% | 0.889±0.129 | +1.23% |
| | FTTransformer-Overlap (ours) | Weak | **0.851±0.138** | / | **0.874±0.130** | / | **0.900±0.127** | / |

**Table A5: AUC-ROC results of ablation studies. Overlap-Gaussian refers to the basic method mentioned in Section 3.3.1. Overlap-Arbitrary refers to the basic method of Eq.4. Overlap-Ranking isolates the ranking loss in Eq.5. Overlap-Combined corresponds to the combined loss form of both Overlap-Arbitrary and Overlap-Ranking as illustrated in Eq.5. Overlap-Proposed refers to the final solution in this paper.**

| Method | $\gamma_l = 5\%$ | | | | | $\gamma_l = 10\%$ | | | | | $\gamma_l = 20\%$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | VAE | MLP | AE | ResNet | FTT | VAE | MLP | AE | ResNet | FTT | VAE | MLP | AE | ResNet | FTT |
| Overlap-Gaussian | 0.539 | / | / | / | / | 0.540 | / | / | / | / | 0.541 | / | / | / | / |
| Overlap-Arbitrary | / | 0.496 | 0.531 | 0.493 | 0.521 | / | 0.498 | 0.543 | 0.534 | 0.482 | / | 0.502 | 0.516 | 0.541 | 0.483 |
| Overlap-Ranking | / | 0.810 | 0.822 | 0.807 | **0.862** | / | 0.845 | 0.857 | 0.855 | **0.888** | / | 0.873 | 0.874 | 0.890 | **0.906** |
| Overlap-Combined | / | 0.843 | 0.854 | 0.820 | 0.610 | / | **0.881** | 0.885 | 0.874 | 0.602 | / | **0.898** | **0.901** | **0.908** | 0.589 |
| Overlap-Proposed | / | **0.847** | **0.862** | **0.836** | 0.851 | / | 0.880 | **0.886** | **0.882** | 0.874 | / | 0.893 | 0.897 | 0.906 | 0.900 |

### Table A6: AUC-ROC results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 5\%$.

| Dataset | Typical | | | | | | MLP | | | AutoEncoder | | | ResNet | | Transformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Iforest | ECOD | Deep SVDD | GAN omaly | Deep SAD | REPEN | DevNet | PReNet | MLP-Overlap | FEAWAD (Sup) | FEAWAD (Weak) | AE-Overlap | ResNet | ResNet-Overlap | FTT | FTT-Overlap |
| ALOI | 0.548 | 0.560 | 0.525 | 0.552 | **0.576** | 0.528 | 0.486 | 0.492 | 0.506 | 0.497 | 0.547 | 0.532 | 0.483 | 0.500 | 0.501 | 0.515 |
| annthyroid | 0.828 | 0.783 | 0.502 | 0.769 | 0.818 | 0.803 | 0.805 | 0.822 | 0.882 | 0.798 | 0.903 | 0.914 | 0.761 | 0.872 | 0.922 | **0.985** |
| Cardiotocography | 0.692 | 0.684 | 0.514 | 0.575 | 0.807 | 0.898 | **0.911** | 0.908 | 0.874 | 0.818 | 0.791 | 0.883 | 0.627 | 0.884 | 0.842 | 0.835 |
| fault | 0.569 | 0.444 | 0.485 | 0.628 | **0.704** | 0.694 | 0.696 | 0.683 | 0.645 | 0.671 | 0.657 | 0.668 | 0.502 | 0.665 | 0.657 | 0.656 |
| http | 0.999 | 0.981 | 0.509 | 0.980 | 0.998 | 0.999 | 0.999 | 0.999 | **1.000** | 0.000 | 0.999 | 0.999 | 0.835 | **1.000** | 0.999 | **1.000** |
| landsat | 0.483 | 0.571 | 0.501 | 0.514 | 0.854 | 0.584 | 0.776 | 0.779 | 0.833 | 0.767 | 0.802 | 0.816 | 0.671 | 0.735 | **0.864** | **0.864** |
| letter | 0.635 | 0.526 | 0.468 | 0.673 | 0.703 | 0.629 | 0.590 | 0.582 | 0.591 | 0.557 | 0.562 | 0.566 | **0.752** | 0.634 | 0.513 | 0.560 |
| magic.gamma | 0.732 | 0.648 | 0.498 | 0.580 | 0.819 | 0.799 | 0.827 | 0.832 | 0.840 | 0.646 | 0.818 | 0.842 | 0.695 | 0.824 | 0.812 | **0.843** |
| mammography | 0.861 | 0.909 | 0.520 | 0.781 | 0.915 | 0.919 | 0.924 | 0.919 | 0.912 | 0.853 | 0.919 | **0.931** | 0.748 | 0.867 | 0.900 | 0.921 |
| mnist | 0.803 | 0.846 | 0.523 | 0.705 | 0.862 | 0.917 | **0.949** | 0.925 | 0.823 | 0.935 | 0.869 | 0.925 | 0.586 | 0.720 | 0.916 | 0.870 |
| musk | 0.999 | 0.952 | 0.550 | 0.781 | 0.917 | 0.915 | **1.000** | **1.000** | **1.000** | 0.998 | 0.927 | **1.000** | 0.426 | **1.000** | 0.999 | 0.989 |
| optdigits | 0.674 | 0.612 | 0.522 | 0.384 | 0.934 | 0.986 | **1.000** | 0.999 | 0.995 | 0.966 | 0.988 | 0.999 | 0.654 | 0.994 | 0.977 | 0.916 |
| PageBlocks | 0.894 | 0.913 | 0.522 | 0.654 | **0.934** | 0.912 | 0.864 | 0.883 | 0.896 | 0.785 | 0.895 | 0.890 | 0.714 | 0.910 | 0.842 | 0.888 |
| pendigits | 0.955 | 0.910 | 0.485 | 0.707 | 0.965 | 0.997 | 0.996 | 0.993 | 0.995 | 0.958 | 0.997 | **0.999** | 0.682 | 0.994 | 0.989 | 0.984 |
| satellite | 0.699 | 0.750 | 0.503 | 0.722 | 0.883 | 0.807 | 0.853 | 0.852 | 0.852 | 0.766 | 0.834 | **0.910** | 0.740 | 0.907 | 0.874 | 0.880 |
| satimage-2 | **0.992** | 0.966 | 0.525 | 0.969 | 0.981 | 0.986 | 0.991 | 0.989 | 0.970 | 0.921 | 0.967 | 0.988 | 0.367 | 0.973 | 0.942 | 0.932 |
| shuttle | **0.996** | 0.995 | 0.508 | 0.744 | 0.990 | 0.989 | 0.979 | 0.978 | 0.981 | 0.976 | 0.980 | 0.982 | 0.973 | 0.979 | 0.976 | 0.977 |
| skin | 0.684 | 0.391 | 0.500 | 0.542 | 0.995 | 0.919 | 0.951 | 0.954 | 0.982 | **0.999** | 0.978 | 0.987 | 0.998 | 0.994 | 0.993 | 0.965 |
| SpamBase | 0.633 | 0.660 | 0.500 | 0.534 | 0.690 | 0.838 | 0.902 | 0.909 | 0.876 | 0.748 | 0.768 | 0.841 | 0.598 | 0.769 | 0.870 | **0.917** |
| speech | 0.498 | 0.510 | 0.549 | 0.481 | 0.531 | 0.582 | 0.604 | 0.631 | 0.589 | 0.587 | 0.475 | 0.624 | 0.532 | 0.581 | 0.551 | **0.648** |
| thyroid | 0.981 | 0.979 | 0.521 | 0.919 | 0.941 | 0.990 | 0.994 | **0.995** | 0.981 | 0.863 | 0.818 | 0.988 | 0.387 | 0.954 | 0.980 | 0.968 |
| vowels | 0.765 | 0.440 | 0.407 | 0.792 | 0.767 | 0.812 | 0.847 | 0.891 | 0.860 | 0.777 | 0.665 | **0.916** | 0.646 | 0.787 | 0.735 | 0.823 |
| Waveform | 0.693 | 0.723 | 0.495 | 0.545 | 0.691 | 0.763 | 0.806 | 0.809 | 0.764 | **0.882** | 0.626 | 0.780 | 0.563 | 0.667 | 0.772 | 0.814 |
| Wilt | 0.427 | 0.395 | 0.502 | 0.388 | 0.799 | 0.529 | 0.689 | 0.695 | 0.922 | 0.894 | 0.819 | 0.945 | 0.795 | **0.950** | 0.658 | 0.910 |
| yeast | 0.382 | 0.387 | 0.477 | 0.449 | 0.512 | 0.467 | 0.625 | 0.627 | 0.611 | 0.625 | 0.610 | 0.619 | 0.530 | **0.642** | 0.583 | 0.627 |

### Table A7: AUC-ROC results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 10\%$.

| Dataset | Typical | | | | | | MLP | | | AutoEncoder | | | ResNet | | Transformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Iforest | ECOD | Deep SVDD | GAN omaly | Deep SAD | REPEN | DevNet | PReNet | MLP-Overlap | FEAWAD (Sup) | FEAWAD (Weak) | AE-Overlap | ResNet | ResNet-Overlap | FTT | FTT-Overlap |
| ALOI | 0.548 | 0.560 | 0.525 | 0.556 | **0.574** | 0.546 | 0.510 | 0.514 | 0.523 | 0.487 | 0.531 | 0.519 | 0.476 | 0.496 | 0.506 | 0.524 |
| annthyroid | 0.828 | 0.783 | 0.502 | 0.733 | 0.884 | 0.824 | 0.826 | 0.834 | 0.939 | 0.880 | 0.897 | 0.968 | 0.905 | 0.932 | **0.990** | 0.990 |
| Cardiotocography | 0.692 | 0.684 | 0.514 | 0.578 | 0.868 | 0.916 | **0.931** | **0.931** | 0.927 | 0.849 | 0.835 | 0.920 | 0.689 | 0.928 | 0.893 | 0.885 |
| fault | 0.569 | 0.444 | 0.485 | 0.631 | **0.728** | 0.720 | 0.724 | 0.719 | 0.695 | 0.692 | 0.674 | 0.695 | 0.570 | 0.721 | 0.694 | 0.699 |
| http | 0.999 | 0.981 | 0.509 | 0.785 | 0.999 | **1.000** | 0.999 | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | 0.829 | **1.000** | **1.000** | **1.000** |
| landsat | 0.483 | 0.571 | 0.501 | 0.522 | **0.897** | 0.561 | 0.779 | 0.789 | 0.878 | 0.805 | 0.805 | 0.839 | 0.743 | 0.836 | 0.891 | 0.891 |
| letter | 0.635 | 0.526 | 0.468 | 0.673 | 0.723 | **0.753** | 0.699 | 0.713 | 0.694 | 0.699 | 0.613 | 0.656 | 0.749 | 0.720 | 0.612 | 0.636 |
| magic.gamma | 0.732 | 0.648 | 0.498 | 0.611 | 0.847 | 0.809 | 0.827 | 0.833 | 0.870 | 0.688 | 0.841 | 0.874 | 0.755 | 0.871 | 0.834 | **0.879** |
| mammography | 0.861 | 0.909 | 0.520 | 0.774 | 0.907 | 0.925 | 0.926 | 0.924 | 0.931 | 0.793 | 0.918 | **0.935** | 0.812 | 0.932 | 0.908 | 0.919 |
| mnist | 0.803 | 0.846 | 0.523 | 0.707 | 0.916 | 0.966 | **0.974** | 0.963 | 0.915 | 0.957 | 0.940 | 0.962 | 0.753 | 0.911 | 0.947 | 0.939 |
| musk | 0.999 | 0.952 | 0.550 | 0.879 | 0.968 | 0.971 | **1.000** | **1.000** | **1.000** | 0.997 | **1.000** | **1.000** | 0.654 | **1.000** | 0.996 | 0.995 |
| optdigits | 0.674 | 0.612 | 0.522 | 0.383 | 0.979 | 0.993 | **1.000** | 0.998 | 0.996 | 0.928 | 0.982 | 0.999 | 0.652 | **1.000** | 0.985 | 0.903 |
| PageBlocks | 0.894 | 0.913 | 0.522 | 0.681 | **0.945** | 0.925 | 0.863 | 0.880 | 0.919 | 0.829 | 0.939 | 0.927 | 0.777 | 0.919 | 0.894 | 0.897 |
| pendigits | 0.955 | 0.910 | 0.485 | 0.713 | 0.993 | 0.996 | 0.995 | 0.995 | 0.995 | 0.981 | 0.992 | **0.999** | 0.761 | **0.999** | 0.997 | 0.981 |
| satellite | 0.699 | 0.750 | 0.503 | 0.726 | 0.909 | 0.807 | 0.851 | 0.849 | 0.899 | 0.840 | 0.852 | **0.919** | 0.759 | 0.910 | 0.908 | 0.891 |
| satimage-2 | **0.992** | 0.966 | 0.525 | 0.969 | 0.986 | 0.989 | 0.988 | 0.986 | 0.954 | 0.959 | 0.971 | 0.978 | 0.801 | 0.987 | 0.974 | 0.949 |
| shuttle | **0.996** | 0.995 | 0.508 | 0.650 | 0.992 | 0.987 | 0.979 | 0.980 | 0.982 | 0.983 | 0.985 | 0.982 | 0.981 | 0.979 | 0.981 | 0.979 |
| skin | 0.684 | 0.391 | 0.500 | 0.516 | 0.997 | 0.910 | 0.951 | 0.954 | 0.992 | **0.999** | 0.992 | 0.996 | 0.998 | 0.993 | 0.995 | 0.985 |
| SpamBase | 0.633 | 0.660 | 0.500 | 0.538 | 0.801 | 0.861 | 0.915 | 0.928 | 0.921 | 0.812 | 0.862 | 0.916 | 0.677 | 0.893 | 0.914 | **0.951** |
| speech | 0.498 | 0.510 | 0.549 | 0.481 | 0.538 | 0.541 | 0.659 | 0.689 | 0.616 | 0.668 | 0.512 | 0.647 | 0.548 | 0.614 | 0.610 | **0.698** |
| thyroid | 0.981 | 0.979 | 0.521 | 0.920 | 0.965 | 0.994 | **0.996** | **0.996** | 0.994 | 0.974 | 0.897 | 0.988 | 0.698 | **0.996** | 0.993 | 0.995 |
| vowels | 0.765 | 0.440 | 0.407 | 0.794 | 0.805 | 0.896 | 0.926 | 0.942 | 0.925 | 0.906 | 0.899 | **0.964** | 0.674 | 0.952 | 0.753 | 0.917 |
| Waveform | 0.693 | 0.723 | 0.495 | 0.545 | 0.765 | 0.846 | 0.881 | 0.884 | 0.836 | **0.887** | 0.728 | 0.854 | 0.666 | 0.786 | 0.863 | 0.774 |
| Wilt | 0.427 | 0.395 | 0.502 | 0.391 | 0.918 | 0.597 | 0.691 | 0.698 | 0.944 | 0.950 | 0.876 | 0.969 | 0.890 | **0.988** | 0.685 | 0.942 |
| yeast | 0.382 | 0.387 | 0.477 | 0.448 | 0.576 | 0.460 | 0.638 | 0.644 | 0.648 | 0.663 | 0.649 | 0.639 | 0.596 | **0.682** | 0.655 | 0.642 |

### Table A8: AUC-ROC results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 20\%$.

| Dataset | Typical | | | | | | MLP | | | AutoEncoder | | | ResNet | | Transformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Iforest | ECOD | Deep SVDD | GAN omaly | Deep SAD | REPEN | DevNet | PReNet | MLP-Overlap | FEAWAD (Sup) | FEAWAD (Weak) | AE-Overlap | ResNet | ResNet-Overlap | FTT | FTT-Overlap |
| ALOI | 0.548 | 0.560 | 0.525 | 0.550 | 0.591 | 0.532 | 0.522 | 0.520 | 0.525 | 0.504 | **0.596** | 0.544 | 0.495 | 0.535 | 0.526 | 0.515 |
| annthyroid | 0.828 | 0.783 | 0.502 | 0.781 | 0.930 | 0.825 | 0.825 | 0.835 | 0.951 | 0.955 | 0.941 | 0.970 | 0.957 | 0.958 | **0.992** | 0.989 |
| Cardiotocography | 0.692 | 0.684 | 0.514 | 0.583 | 0.913 | 0.930 | 0.946 | 0.945 | 0.942 | 0.896 | 0.894 | 0.946 | 0.778 | **0.948** | 0.930 | 0.933 |
| fault | 0.569 | 0.444 | 0.485 | 0.633 | **0.749** | 0.743 | 0.733 | 0.738 | 0.698 | 0.669 | 0.718 | 0.673 | 0.642 | 0.745 | 0.713 | 0.717 |
| http | 0.999 | 0.981 | 0.509 | 0.783 | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** |
| landsat | 0.483 | 0.571 | 0.501 | 0.538 | **0.921** | 0.560 | 0.789 | 0.782 | 0.899 | 0.872 | 0.846 | 0.806 | 0.806 | 0.884 | 0.913 | 0.901 |
| letter | 0.635 | 0.526 | 0.468 | 0.675 | 0.751 | **0.806** | 0.748 | 0.767 | 0.739 | 0.717 | 0.677 | 0.683 | 0.739 | 0.779 | 0.705 | 0.720 |
| magic.gamma | 0.732 | 0.648 | 0.498 | 0.662 | 0.874 | 0.812 | 0.827 | 0.831 | 0.874 | 0.746 | 0.828 | 0.881 | 0.803 | **0.889** | 0.858 | 0.860 |
| mammography | 0.861 | 0.909 | 0.520 | 0.755 | 0.934 | 0.932 | 0.927 | 0.925 | 0.945 | 0.856 | 0.934 | **0.948** | 0.864 | 0.946 | 0.923 | 0.934 |
| mnist | 0.803 | 0.846 | 0.523 | 0.705 | 0.959 | **0.984** | **0.984** | 0.983 | 0.969 | 0.967 | 0.965 | 0.982 | 0.846 | 0.980 | 0.968 | 0.979 |
| musk | 0.999 | 0.952 | 0.550 | 0.891 | 0.996 | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | 0.752 | **1.000** | **1.000** | **1.000** |
| optdigits | 0.674 | 0.612 | 0.522 | 0.385 | 0.997 | 0.996 | **1.000** | **1.000** | 0.998 | 0.987 | **1.000** | 0.999 | 0.777 | **1.000** | 0.998 | 0.999 |
| PageBlocks | 0.894 | 0.913 | 0.522 | 0.775 | **0.959** | 0.929 | 0.885 | 0.897 | 0.935 | 0.901 | 0.948 | 0.949 | 0.905 | 0.933 | 0.925 | 0.933 |
| pendigits | 0.955 | 0.910 | 0.485 | 0.718 | 0.999 | 0.996 | 0.997 | 0.997 | 0.998 | 0.993 | 0.996 | **1.000** | 0.923 | 0.999 | 0.997 | 0.993 |
| satellite | 0.699 | 0.750 | 0.503 | 0.739 | 0.932 | 0.806 | 0.856 | 0.849 | 0.903 | 0.867 | 0.866 | 0.928 | 0.836 | **0.937** | 0.935 | 0.914 |
| satimage-2 | 0.992 | 0.966 | 0.525 | 0.971 | 0.992 | 0.992 | **0.993** | **0.993** | 0.969 | 0.959 | 0.977 | 0.986 | 0.912 | 0.988 | 0.969 | 0.974 |
| shuttle | **0.996** | 0.995 | 0.508 | 0.757 | 0.992 | 0.987 | 0.979 | 0.979 | 0.983 | 0.987 | 0.987 | 0.984 | 0.982 | 0.980 | 0.984 | 0.978 |
| skin | 0.684 | 0.391 | 0.500 | 0.499 | 0.998 | 0.903 | 0.951 | 0.955 | 0.995 | 0.999 | 0.983 | 0.998 | 0.999 | **1.000** | 0.996 | 0.976 |
| SpamBase | 0.633 | 0.660 | 0.500 | 0.544 | 0.887 | 0.889 | 0.913 | 0.931 | 0.938 | 0.868 | 0.886 | 0.946 | 0.807 | 0.945 | 0.928 | **0.953** |
| speech | 0.498 | 0.510 | 0.549 | 0.482 | 0.559 | 0.611 | 0.713 | **0.740** | 0.624 | 0.643 | 0.598 | 0.694 | 0.584 | 0.685 | 0.678 | 0.734 |
| thyroid | 0.981 | 0.979 | 0.521 | 0.918 | 0.986 | 0.996 | **0.997** | **0.997** | 0.994 | 0.983 | 0.996 | 0.996 | 0.879 | **0.997** | 0.993 | **0.997** |
| vowels | 0.765 | 0.440 | 0.407 | 0.798 | 0.871 | 0.971 | 0.970 | 0.979 | 0.977 | 0.953 | 0.961 | 0.995 | 0.837 | **0.996** | 0.976 | 0.991 |
| Waveform | 0.693 | 0.723 | 0.495 | 0.546 | 0.811 | 0.898 | **0.909** | 0.903 | 0.880 | 0.899 | 0.765 | 0.885 | 0.700 | 0.839 | 0.885 | 0.885 |
| Wilt | 0.427 | 0.395 | 0.502 | 0.479 | 0.960 | 0.646 | 0.689 | 0.696 | 0.954 | 0.981 | 0.876 | 0.970 | 0.922 | **0.991** | 0.732 | 0.976 |
| yeast | 0.382 | 0.387 | 0.477 | 0.458 | 0.644 | 0.459 | 0.663 | 0.666 | 0.646 | 0.703 | 0.663 | 0.672 | 0.660 | 0.699 | **0.705** | 0.650 |

### Table A9: AUC-PR results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 5\%$.

| Dataset | Typical | | | | | | MLP | | | AutoEncoder | | | ResNet | | Transformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Iforest | ECOD | Deep SVDD | GAN omaly | Deep SAD | REPEN | DevNet | PReNet | MLP-Overlap | FEAWAD (Sup) | FEAWAD (Weak) | AE-Overlap | ResNet | ResNet-Overlap | FTT | FTT-Overlap |
| ALOI | 0.036 | 0.036 | 0.042 | 0.038 | **0.058** | 0.041 | 0.042 | 0.045 | 0.046 | 0.043 | 0.040 | 0.048 | 0.037 | 0.040 | 0.035 | 0.036 |
| annthyroid | 0.336 | 0.260 | 0.078 | 0.359 | 0.402 | 0.400 | 0.425 | 0.458 | 0.541 | 0.482 | 0.609 | 0.602 | 0.363 | 0.575 | 0.658 | **0.799** |
| Cardiotocography | 0.441 | 0.436 | 0.254 | 0.345 | 0.602 | 0.760 | 0.767 | **0.776** | 0.708 | 0.610 | 0.652 | 0.720 | 0.433 | 0.736 | 0.649 | 0.691 |
| fault | 0.418 | 0.317 | 0.355 | 0.477 | 0.531 | 0.543 | 0.557 | **0.565** | 0.514 | 0.501 | 0.527 | 0.525 | 0.425 | 0.521 | 0.495 | 0.538 |
| http | 0.808 | 0.158 | 0.022 | 0.648 | 0.659 | 0.833 | 0.869 | 0.845 | 0.982 | 0.003 | 0.847 | 0.898 | 0.805 | **1.000** | 0.867 | **1.000** |
| landsat | 0.199 | 0.262 | 0.210 | 0.210 | 0.617 | 0.322 | 0.448 | 0.447 | 0.572 | 0.502 | 0.575 | 0.537 | 0.407 | 0.497 | **0.682** | 0.633 |
| letter | 0.098 | 0.076 | 0.165 | 0.141 | 0.176 | 0.148 | 0.158 | 0.159 | 0.201 | 0.134 | 0.109 | 0.137 | **0.342** | 0.169 | 0.118 | 0.126 |
| magic.gamma | 0.648 | 0.551 | 0.357 | 0.449 | 0.709 | 0.730 | 0.725 | 0.755 | 0.771 | 0.535 | 0.732 | 0.774 | 0.548 | 0.750 | 0.703 | **0.792** |
| mammography | 0.195 | 0.414 | 0.055 | 0.127 | 0.510 | 0.592 | **0.610** | 0.603 | 0.427 | 0.418 | 0.545 | 0.546 | 0.372 | 0.469 | 0.528 | 0.589 |
| mnist | 0.286 | 0.329 | 0.127 | 0.196 | 0.544 | 0.717 | **0.784** | 0.738 | 0.617 | 0.670 | 0.643 | 0.760 | 0.286 | 0.496 | 0.644 | 0.554 |
| musk | 0.970 | 0.343 | 0.105 | 0.464 | 0.597 | 0.653 | **1.000** | **1.000** | **1.000** | 0.961 | 0.895 | **1.000** | 0.246 | **1.000** | 0.974 | 0.951 |
| optdigits | 0.047 | 0.035 | 0.032 | 0.028 | 0.510 | 0.951 | **0.991** | 0.989 | 0.966 | 0.708 | 0.961 | 0.987 | 0.260 | 0.960 | 0.811 | 0.625 |
| PageBlocks | 0.469 | 0.517 | 0.146 | 0.346 | 0.714 | 0.678 | 0.654 | 0.679 | 0.642 | 0.539 | **0.719** | 0.680 | 0.456 | 0.696 | 0.602 | 0.638 |
| pendigits | 0.283 | 0.216 | 0.032 | 0.185 | 0.716 | 0.955 | 0.928 | 0.922 | 0.952 | 0.760 | 0.954 | 0.980 | 0.367 | 0.970 | 0.890 | 0.908 |
| satellite | 0.662 | 0.662 | 0.323 | 0.666 | 0.782 | 0.795 | 0.828 | 0.825 | 0.768 | 0.663 | 0.763 | 0.861 | 0.628 | **0.869** | 0.805 | 0.838 |
| satimage-2 | 0.913 | 0.629 | 0.034 | 0.523 | 0.671 | 0.899 | 0.908 | 0.912 | 0.918 | 0.560 | 0.912 | **0.930** | 0.161 | 0.880 | 0.885 | 0.790 |
| shuttle | **0.975** | 0.957 | 0.081 | 0.464 | 0.954 | 0.974 | 0.968 | 0.967 | 0.962 | 0.956 | 0.967 | 0.966 | 0.959 | 0.966 | 0.953 | 0.960 |
| skin | 0.257 | 0.156 | 0.204 | 0.223 | 0.950 | 0.555 | 0.660 | 0.673 | 0.866 | **0.994** | 0.832 | 0.905 | 0.984 | 0.950 | 0.963 | 0.793 |
| SpamBase | 0.496 | 0.528 | 0.400 | 0.409 | 0.584 | 0.774 | 0.846 | 0.863 | 0.836 | 0.657 | 0.754 | 0.806 | 0.550 | 0.748 | 0.813 | **0.889** |
| speech | 0.021 | 0.019 | 0.049 | 0.017 | 0.024 | 0.045 | 0.052 | 0.057 | 0.042 | 0.044 | 0.027 | 0.060 | 0.040 | 0.046 | 0.037 | 0.054 |
| thyroid | 0.574 | 0.526 | 0.068 | 0.466 | 0.464 | 0.753 | 0.880 | **0.888** | 0.814 | 0.487 | 0.581 | 0.852 | 0.108 | 0.777 | 0.737 | 0.814 |
| vowels | 0.193 | 0.039 | 0.105 | 0.244 | 0.150 | 0.345 | 0.384 | 0.431 | 0.438 | 0.289 | 0.283 | **0.602** | 0.336 | 0.354 | 0.336 | 0.386 |
| Waveform | 0.063 | 0.064 | 0.032 | 0.043 | 0.180 | 0.121 | 0.135 | 0.177 | 0.182 | **0.216** | 0.156 | 0.181 | 0.141 | 0.154 | 0.182 | 0.190 |
| Wilt | 0.043 | 0.042 | 0.054 | 0.045 | 0.188 | 0.065 | 0.087 | 0.089 | 0.381 | 0.545 | 0.390 | 0.494 | 0.380 | 0.582 | 0.082 | **0.633** |
| yeast | 0.293 | 0.305 | 0.333 | 0.314 | 0.348 | 0.338 | 0.437 | 0.439 | 0.436 | 0.445 | 0.429 | 0.438 | 0.391 | **0.467** | 0.409 | 0.447 |

### Table A10: AUC-PR results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 10\%$.

| Dataset | Typical | | | | | | MLP | | | AutoEncoder | | | ResNet | | Transformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Iforest | ECOD | Deep SVDD | GAN omaly | Deep SAD | REPEN | DevNet | PReNet | MLP-Overlap | FEAWAD (Sup) | FEAWAD (Weak) | AE-Overlap | ResNet | ResNet-Overlap | FTT | FTT-Overlap |
| ALOI | 0.036 | 0.036 | 0.042 | 0.039 | **0.068** | 0.052 | 0.047 | 0.052 | 0.049 | 0.044 | 0.041 | 0.046 | 0.036 | 0.037 | 0.037 | 0.042 |
| annthyroid | 0.336 | 0.260 | 0.078 | 0.327 | 0.506 | 0.437 | 0.459 | 0.476 | 0.660 | 0.615 | 0.574 | 0.729 | 0.591 | 0.674 | **0.835** | 0.833 |
| Cardiotocography | 0.441 | 0.436 | 0.254 | 0.347 | 0.708 | 0.799 | 0.805 | **0.817** | 0.784 | 0.676 | 0.710 | 0.779 | 0.501 | 0.799 | 0.695 | 0.755 |
| fault | 0.418 | 0.317 | 0.355 | 0.480 | 0.560 | 0.577 | 0.586 | **0.596** | 0.542 | 0.521 | 0.543 | 0.537 | 0.490 | 0.575 | 0.538 | 0.567 |
| http | 0.808 | 0.158 | 0.022 | 0.451 | 0.829 | 0.928 | 0.891 | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | 0.801 | **1.000** | **1.000** | **1.000** |
| landsat | 0.199 | 0.262 | 0.210 | 0.215 | 0.704 | 0.357 | 0.454 | 0.486 | 0.632 | 0.568 | 0.516 | 0.597 | 0.474 | 0.614 | **0.719** | 0.708 |
| letter | 0.098 | 0.076 | 0.165 | 0.142 | 0.183 | 0.209 | 0.153 | 0.170 | 0.225 | 0.167 | 0.140 | 0.144 | **0.305** | 0.190 | 0.139 | 0.142 |
| magic.gamma | 0.648 | 0.551 | 0.357 | 0.473 | 0.754 | 0.748 | 0.718 | 0.748 | 0.817 | 0.594 | 0.758 | 0.822 | 0.622 | 0.819 | 0.735 | **0.839** |
| mammography | 0.195 | 0.414 | 0.055 | 0.127 | 0.562 | 0.606 | **0.621** | 0.606 | 0.515 | 0.370 | 0.520 | 0.566 | 0.475 | 0.573 | 0.565 | 0.544 |
| mnist | 0.286 | 0.329 | 0.127 | 0.198 | 0.678 | **0.847** | 0.846 | 0.820 | 0.752 | 0.737 | 0.737 | 0.844 | 0.443 | 0.781 | 0.724 | 0.799 |
| musk | 0.970 | 0.343 | 0.105 | 0.733 | 0.835 | 0.886 | **1.000** | **1.000** | **1.000** | 0.972 | **1.000** | **1.000** | 0.405 | **1.000** | 0.983 | 0.988 |
| optdigits | 0.047 | 0.035 | 0.032 | 0.027 | 0.811 | 0.973 | 0.994 | 0.991 | 0.979 | 0.713 | 0.961 | 0.992 | 0.323 | **0.995** | 0.885 | 0.746 |
| PageBlocks | 0.469 | 0.517 | 0.146 | 0.349 | 0.747 | 0.719 | 0.645 | 0.662 | 0.666 | 0.616 | **0.774** | 0.695 | 0.526 | 0.695 | 0.689 | 0.711 |
| pendigits | 0.283 | 0.216 | 0.032 | 0.186 | 0.893 | 0.963 | 0.914 | 0.914 | 0.955 | 0.881 | 0.946 | **0.985** | 0.427 | 0.980 | 0.946 | 0.865 |
| satellite | 0.662 | 0.662 | 0.323 | 0.653 | 0.829 | 0.789 | 0.830 | 0.828 | 0.846 | 0.749 | 0.772 | **0.877** | 0.636 | 0.872 | 0.829 | 0.851 |
| satimage-2 | 0.913 | 0.629 | 0.034 | 0.527 | 0.868 | 0.910 | 0.897 | 0.886 | 0.911 | 0.851 | 0.907 | **0.919** | 0.525 | 0.916 | 0.873 | 0.877 |
| shuttle | **0.975** | 0.957 | 0.081 | 0.363 | 0.964 | 0.971 | 0.968 | 0.967 | 0.966 | 0.971 | 0.970 | 0.971 | 0.967 | 0.968 | 0.961 | 0.965 |
| skin | 0.257 | 0.156 | 0.204 | 0.212 | 0.971 | 0.532 | 0.658 | 0.675 | 0.937 | **0.989** | 0.945 | 0.959 | 0.982 | 0.944 | 0.970 | 0.904 |
| SpamBase | 0.496 | 0.528 | 0.400 | 0.411 | 0.701 | 0.802 | 0.857 | 0.879 | 0.884 | 0.745 | 0.834 | 0.886 | 0.635 | 0.868 | 0.862 | **0.934** |
| speech | 0.021 | 0.019 | 0.049 | 0.017 | 0.025 | 0.056 | 0.065 | 0.068 | **0.085** | 0.066 | 0.044 | 0.065 | 0.064 | 0.067 | 0.061 | 0.072 |
| thyroid | 0.574 | 0.526 | 0.068 | 0.475 | 0.606 | 0.870 | **0.903** | 0.883 | 0.870 | 0.750 | 0.810 | 0.882 | 0.368 | 0.895 | 0.840 | 0.893 |
| vowels | 0.193 | 0.039 | 0.105 | 0.247 | 0.192 | 0.477 | 0.646 | 0.705 | 0.676 | 0.506 | 0.501 | **0.770** | 0.410 | 0.676 | 0.418 | 0.711 |
| Waveform | 0.063 | 0.064 | 0.032 | 0.043 | **0.263** | 0.146 | 0.160 | 0.189 | 0.221 | 0.245 | 0.195 | 0.248 | 0.137 | 0.218 | 0.232 | 0.220 |
| Wilt | 0.043 | 0.042 | 0.054 | 0.046 | 0.387 | 0.077 | 0.087 | 0.089 | 0.436 | 0.663 | 0.462 | 0.611 | 0.492 | **0.831** | 0.090 | 0.752 |
| yeast | 0.293 | 0.305 | 0.333 | 0.313 | 0.392 | 0.349 | 0.434 | 0.440 | 0.443 | **0.489** | 0.462 | 0.439 | 0.429 | 0.486 | 0.472 | 0.443 |

### Table A11: AUC-PR results of model comparison on 25 real-world datasets w.r.t. $\gamma_l = 20\%$.

| Dataset | Typical | | | | | | MLP | | | AutoEncoder | | | ResNet | | Transformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Iforest | ECOD | Deep SVDD | GAN omaly | Deep SAD | REPEN | DevNet | PReNet | MLP-Overlap | FEAWAD (Sup) | FEAWAD (Weak) | AE-Overlap | ResNet | ResNet-Overlap | FTT | FTT-Overlap |
| ALOI | 0.036 | 0.036 | 0.042 | 0.038 | **0.069** | 0.046 | 0.047 | 0.045 | 0.045 | 0.051 | 0.047 | 0.044 | 0.040 | 0.052 | 0.041 | 0.039 |
| annthyroid | 0.336 | 0.260 | 0.078 | 0.372 | 0.624 | 0.441 | 0.457 | 0.478 | 0.677 | 0.703 | 0.622 | 0.739 | 0.704 | 0.748 | **0.858** | 0.818 |
| Cardiotocography | 0.441 | 0.436 | 0.254 | 0.354 | 0.793 | 0.830 | **0.838** | 0.837 | 0.826 | 0.743 | 0.768 | 0.835 | 0.600 | **0.838** | 0.794 | 0.828 |
| fault | 0.418 | 0.317 | 0.355 | 0.485 | 0.582 | **0.610** | 0.581 | 0.604 | 0.510 | 0.517 | 0.570 | 0.520 | 0.536 | 0.592 | 0.548 | 0.565 |
| http | 0.808 | 0.158 | 0.022 | 0.451 | 0.891 | 0.928 | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** |
| landsat | 0.199 | 0.262 | 0.210 | 0.228 | **0.763** | 0.389 | 0.528 | 0.487 | 0.688 | 0.695 | 0.563 | 0.523 | 0.587 | 0.680 | 0.753 | 0.745 |
| letter | 0.098 | 0.076 | 0.165 | 0.143 | 0.197 | 0.276 | 0.216 | 0.260 | 0.239 | 0.206 | 0.195 | 0.175 | **0.344** | 0.268 | 0.197 | 0.211 |
| magic.gamma | 0.648 | 0.551 | 0.357 | 0.518 | 0.799 | 0.746 | 0.720 | 0.743 | 0.817 | 0.674 | 0.726 | 0.826 | 0.684 | **0.850** | 0.767 | 0.819 |
| mammography | 0.195 | 0.414 | 0.055 | 0.119 | 0.600 | 0.617 | 0.614 | 0.613 | 0.513 | 0.477 | 0.618 | 0.600 | 0.510 | 0.602 | 0.601 | **0.629** |
| mnist | 0.286 | 0.329 | 0.127 | 0.195 | 0.810 | **0.908** | 0.871 | 0.888 | 0.850 | 0.808 | 0.786 | 0.883 | 0.598 | 0.902 | 0.798 | 0.872 |
| musk | 0.970 | 0.343 | 0.105 | 0.752 | 0.981 | 0.994 | **1.000** | **1.000** | **1.000** | 0.999 | **1.000** | **1.000** | 0.592 | **1.000** | **1.000** | 0.999 |
| optdigits | 0.047 | 0.035 | 0.032 | 0.028 | 0.961 | 0.989 | 0.996 | 0.994 | 0.987 | 0.887 | 0.993 | 0.992 | 0.517 | **0.999** | 0.967 | 0.982 |
| PageBlocks | 0.469 | 0.517 | 0.146 | 0.395 | **0.787** | 0.724 | 0.672 | 0.692 | 0.696 | 0.725 | 0.785 | 0.752 | 0.720 | 0.710 | 0.708 | 0.738 |
| pendigits | 0.283 | 0.216 | 0.032 | 0.188 | 0.984 | 0.973 | 0.933 | 0.940 | 0.982 | 0.919 | 0.965 | **0.990** | 0.733 | 0.986 | 0.969 | 0.958 |
| satellite | 0.662 | 0.662 | 0.323 | 0.646 | 0.874 | 0.790 | 0.832 | 0.829 | 0.854 | 0.794 | 0.801 | 0.888 | 0.741 | **0.903** | 0.870 | 0.882 |
| satimage-2 | 0.913 | 0.629 | 0.034 | 0.535 | 0.906 | 0.923 | **0.928** | 0.919 | 0.926 | 0.877 | 0.926 | 0.920 | 0.752 | 0.927 | 0.907 | 0.911 |
| shuttle | **0.975** | 0.957 | 0.081 | 0.415 | 0.967 | 0.972 | 0.967 | 0.967 | 0.966 | 0.973 | 0.971 | 0.972 | 0.968 | 0.969 | 0.964 | 0.969 |
| skin | 0.257 | 0.156 | 0.204 | 0.202 | 0.981 | 0.513 | 0.658 | 0.679 | 0.954 | 0.988 | 0.871 | 0.988 | 0.990 | **0.998** | 0.974 | 0.857 |
| SpamBase | 0.496 | 0.528 | 0.400 | 0.416 | 0.811 | 0.838 | 0.858 | 0.885 | 0.905 | 0.814 | 0.854 | 0.923 | 0.759 | 0.928 | 0.879 | **0.938** |
| speech | 0.021 | 0.019 | 0.049 | 0.017 | 0.027 | **0.160** | 0.101 | 0.117 | 0.084 | 0.075 | 0.074 | 0.070 | 0.119 | 0.108 | 0.063 | 0.104 |
| thyroid | 0.574 | 0.526 | 0.068 | 0.475 | 0.785 | 0.866 | 0.911 | 0.898 | 0.859 | 0.751 | 0.871 | 0.897 | 0.574 | 0.891 | **0.916** | 0.914 |
| vowels | 0.193 | 0.039 | 0.105 | 0.253 | 0.341 | 0.833 | 0.807 | 0.838 | 0.825 | 0.654 | 0.803 | 0.921 | 0.621 | 0.925 | 0.757 | **0.941** |
| Waveform | 0.063 | 0.064 | 0.032 | 0.043 | **0.335** | 0.185 | 0.217 | 0.227 | 0.282 | 0.286 | 0.224 | 0.310 | 0.192 | 0.282 | 0.328 | 0.265 |
| Wilt | 0.043 | 0.042 | 0.054 | 0.053 | 0.573 | 0.087 | 0.086 | 0.088 | 0.476 | 0.799 | 0.537 | 0.579 | 0.590 | **0.874** | 0.110 | 0.822 |
| yeast | 0.293 | 0.305 | 0.333 | 0.318 | 0.445 | 0.349 | 0.462 | 0.467 | 0.438 | **0.535** | 0.477 | 0.476 | 0.493 | 0.509 | 0.514 | 0.453 |

## .5 Additional Results of AD Loss Function Exploration

In addition to the main paper that demonstrates the embedding variations on the vowels dataset in Section 4.3.1, here we provide another example of the skin dataset, as shown in Figure A1. Compared to the other loss functions, our proposed Overlap loss better retains the ringlike shape in the embedding of input feature while achieving satisfactory detection performance.

We follow [21, 58] to generate the following four types of synthetic anomalies, which are further used to evaluate different AD loss functions. The AUC-ROC results of loss function comparison on different types of anomalies also indicate that our proposed Overlap loss significantly outperforms other counterparts, as is shown in Table A12.

- **Local anomalies** refer to the anomalies deviant from their local neighborhoods [9]. GMM procedure [43, 58] is used to generate synthetic normal samples, and then scale the covariance matrix $\hat{\Sigma} = \alpha \hat{\Sigma}$ by a scaling parameter $\alpha = 5$ to generate local anomalies.
- **Global anomalies** are generated from a uniform distribution $\text{Unif}\left(\alpha \cdot \min\left(x^k\right), \alpha \cdot \max\left(x^k\right)\right)$, where the boundaries are defined as the *min* and *max* of an input feature, e.g., $k$-th feature $x^k$, and $\alpha = 1.1$ controls the outlyingness of anomalies.
- **Dependency anomalies** refer to the samples that do not follow the dependency structure that normal data follows [42], i.e., the input features of dependency anomalies are assumed to be independent of each other. Vine Copula [1] method is applied to model the dependency structure of original data, where the probability density function of generated anomalies is set to complete independence by removing the modeled dependency (see [42]). KDE method estimates the probability density function of features and generates normal samples.

- **Clustered anomalies**, also known as group anomalies [32], exhibit similar characteristics [15, 36]. We scale the mean feature vector of normal samples by $\alpha = 5$, i.e., $\hat{\mu} = \alpha \hat{\mu}$, where $\alpha$ controls the distance between anomaly clusters and the normals, and use the scaled GMM to generate anomalies.

**Table A12: Loss function comparison on different types of anomalies generated based on the 25 real-world datasets.**

| Loss | Local | Global | Clustered | Dependency |
|------|-------|--------|-----------|------------|
| Minus | 0.629 | 0.936 | 0.996 | 0.738 |
| Inverse | 0.547 | 0.823 | 0.937 | 0.570 |
| Hinge | 0.607 | 0.938 | 0.997 | 0.761 |
| Deviation | 0.588 | 0.959 | 0.990 | 0.652 |
| Ordinal | 0.604 | 0.954 | 0.994 | 0.687 |
| Overlap | **0.742** | **0.981** | **0.998** | **0.847** |

We further investigate two case studies by generating visualized two-dimensional synthetic samples of the above local and clustered anomalies, as shown in Figure A2. The anomaly ratios of these two datasets are set to 5%. The results indicate that all the compared loss functions can correctly detect anomalies for the two-dimensional clustered anomalies (with 1.000 AUC-ROC and AUC-PR). This result can be expected since few labeled clustered anomalies can already represent similar behaviors of the entire clustered anomalies. For the local anomalies, however, we observe most of the compared loss functions perform poorly. In contrast, Overlap loss achieves better detection performance, and successfully learns a suitable decision boundary (see Figure A2l), where the learned decision boundary fits well with the local anomalies that are often overlapped or surrounded by the normal samples.

(a) t-SNE [60] plots of the input feature of skin dataset.

(b) Minus

(c) Inverse

(d) Hinge

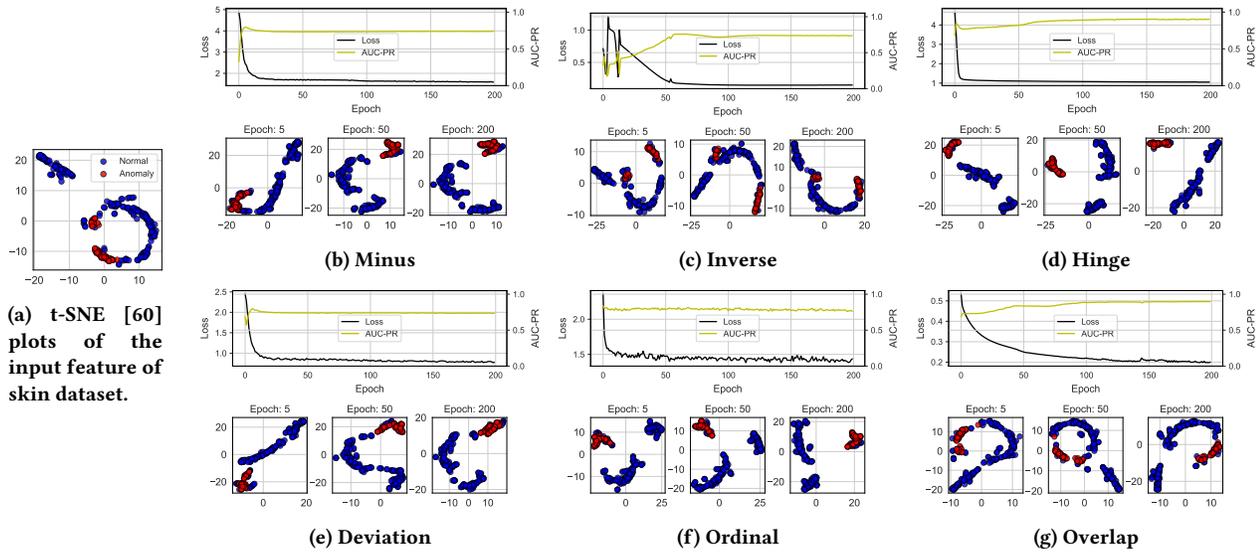(e) Deviation

(f) Ordinal

(g) Overlap

**Figure A1: Training loss along with the AUC-PR performance on testing set of different loss function based AD models, where the skin dataset is specified for comparison. The transformed embeddings of the input feature are demonstrated, which corresponds to 5, 50, and 200 training epochs, respectively.**



Clustered Anomaly

(a) Minus

(b) Inverse

(c) Hinge

(d) Deviation

(e) Ordinal

(f) Overlap

Local Anomaly

(g) Minus

(h) Inverse

(i) Hinge

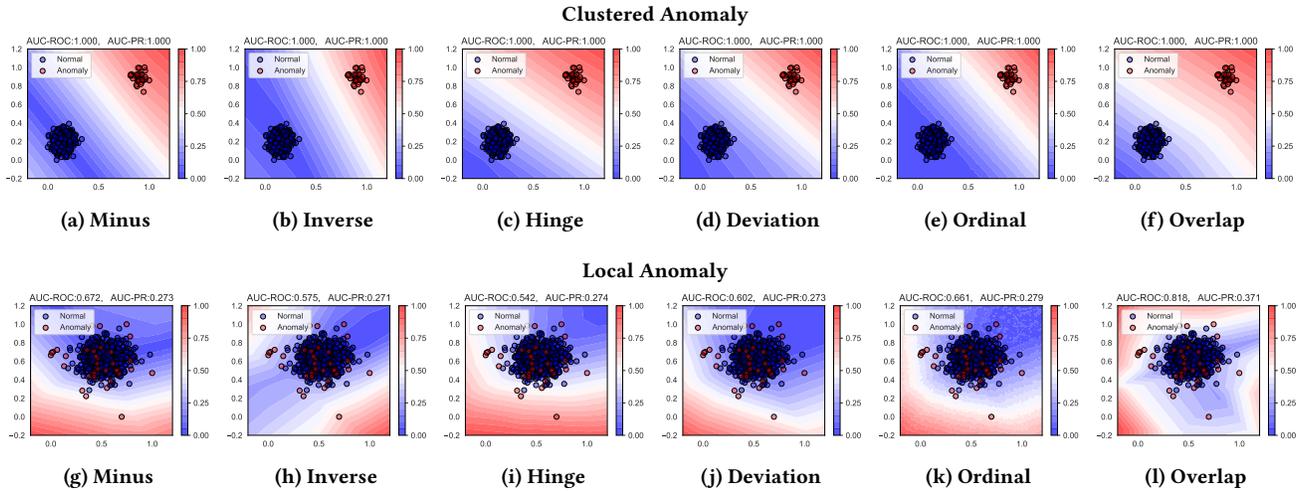(j) Deviation

(k) Ordinal

(l) Overlap

**Figure A2: Decision boundaries of different loss functions on the local anomalies. The output anomaly scores are normalized to $[0, 1]$ for comparison. Both AUC-ROC and AUC-PR performances are displayed in the title above each subfigure.**