
ADBench: Anomaly Detection Benchmark

Songqiao Han^{1,*}, Xiyang Hu^{2,*}, Hailiang Huang^{1,*}, Mingqi Jiang^{1,*}, Yue Zhao^{2,*}

¹ Shanghai University of Finance and Economics ² Carnegie Mellon University
{han.songqiao, hlhuang}@shufe.edu.cn, {2020310191}@live.sufe.edu.cn,
{xiyanghu, zhaoy}@cmu.edu

Abstract

Given a long list of anomaly detection algorithms developed in the last few decades, how do they perform with regard to (i) varying levels of supervision, (ii) different types of anomalies, and (iii) noisy and corrupted data? In this work, we answer these key questions by conducting (to our best knowledge) the most comprehensive anomaly detection benchmark with 30 algorithms on 55 benchmark datasets, named ADBench. Our extensive experiments (93,654 in total) identify meaningful insights into the role of supervision and anomaly types, and unlock future directions for researchers in algorithm selection and design. With ADBench, researchers can easily conduct comprehensive and fair evaluations for newly proposed methods on the datasets (including our contributed ones from natural language and computer vision domains) against the existing baselines. To foster accessibility and reproducibility, we fully open-source ADBench and the corresponding results.

1 Introduction

Anomaly detection (AD), which is also known as outlier detection, is a key machine learning (ML) task with numerous applications, including anti-money laundering [57], rare disease detection [116], social media analysis [110, 114], and intrusion detection [54]. AD algorithms aim to identify data instances that deviate significantly from the majority of data objects [35, 82, 87, 96], and numerous methods have been developed in the last few decades [2, 53, 64, 65, 76, 93, 104, 118]. Among them, majority are designed for tabular data (i.e., no time dependency and graph structure). Thus, we focuses on the *tabular* AD algorithms and datasets in this work.

Although there are already some benchmark and evaluation works for tabular AD [14, 23, 25, 30, 99], they generally have the limitations as follows: (i) primary emphasis on unsupervised methods only without including emerging (semi-)supervised AD methods; (ii) limited analysis of the algorithm performance w.r.t. anomaly types; (iii) the lack of analysis on model robustness (e.g., noisy labels and irrelevant features); (iv) the absence of using statistical tests for algorithm comparison; and (v) no coverage of more complex NLP and CV datasets, which have attracted extensive attention nowadays.

To address these limitations, we design (to our best knowledge) the most comprehensive tabular anomaly detection benchmark called ADBench. By analyzing both research needs and deployment requirements in industry, we design the experiments with three major angles in anomaly detection (see §3.3): (i) the availability of supervision (e.g., ground truth labels) by including 14 unsupervised, 7 semi-supervised, and 9 supervised methods; (ii) algorithm performance under different types of anomalies by simulating the environments with 4 types of anomalies; and (iii) algorithm robustness and stability under 3 settings of data corruptions. Fig. 1 provides an overview of ADBench.

*All authors contribute equally and are listed alphabetically.

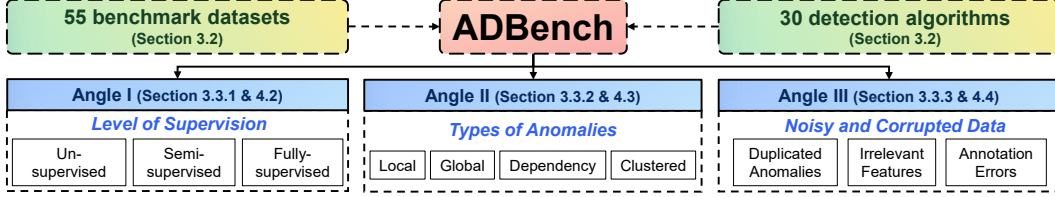


Figure 1: The primary design of the proposed ADBench driven by research and application needs.

Key takeaways: Through extensive experiments, we find (i) surprisingly none of the benchmarked unsupervised algorithms is statistically better than others, emphasizing the importance of algorithm selection; (ii) with merely 1% labeled anomalies, most semi-supervised methods can outperform the best unsupervised method, justifying the importance of supervision; (iii) in controlled environments, we observe that best unsupervised methods for specific types of anomalies are even better than semi- and fully-supervised methods, revealing the necessity of understanding data characteristics; (iv) semi-supervised methods show potential in achieving robustness in noisy and corrupted data, possibly due to their efficiency in using labels and feature selection. See §4 for additional results and insights.

We summarize the primary contributions of ADBench as below:

1. **The most comprehensive AD benchmark.** ADBench examines 30 detection algorithms’ performance on 55 benchmark datasets (of which 47 are existing ones and 8 are created by us).
2. **Research and application-driven benchmark angles.** By analyzing the needs of research and real-world applications, we focus on three key comparison angles, namely the availability of supervision, the anomaly types, and the algorithm robustness under noise and data corruption.
3. **Insights and future directions for researchers and practitioners.** Based on the extensive results, we show the necessity of algorithm selection, and the value of supervision and prior knowledge.
4. **Fair and accessible AD evaluation.** We open-source ADBench with BSD-2 License at <https://github.com/Minqi824/ADbench>, for benchmarking on newly proposed methods.

2 Related Work

2.1 Anomaly Detection Algorithms

Unsupervised Methods by Assuming Anomaly Data Distributions. *Unsupervised AD methods are proposed with different assumptions of data distribution* [2], e.g., anomalies locate in low-density regions, and their performance often depends on the agreement between the input data and the algorithm assumption(s). Many unsupervised methods have been proposed in the last few decades [2, 10, 76, 90, 118], which can be roughly categorized into shallow and deep (neural network) methods. The former often carries better interpretability, while the latter handles large, high-dimensional data better. Please see Appx. §A.1, recent book [2], and surveys [76, 90] for additional information.

Supervised Methods by Treating Anomaly Detection as Binary Classification. *With the accessibility of full ground truth labels (which is rare), supervised classifiers may identify known anomalies at the risk of missing unknown anomalies.* Arguably, there are no specialized supervised anomaly detection algorithms, and people often use existing classifiers for this purpose [2] such as Random Forest [11] and neural networks [55]. One known risk of supervised methods is that ground truth labels may be not necessarily accurate enough to capture all types of anomalies during annotation, and these methods are therefore limited to detect unknown types of anomalies [2]. Recent machine learning books [3, 31] and scikit-learn [79] may serve as good sources of supervised ML methods.

Semi-supervised Methods with Efficient Use of Labels. *Semi-supervised AD algorithms are designed to capitalize the supervision from partial labels, while keeping the ability of detecting unseen types of anomalies.* To this end, some recent studies investigate to efficiently use partially labeled data for improving detection performance, and leverage the unlabeled data to facilitate representation learning. For instance, some semi-supervised models are trained only on normal samples, and detect anomalies that deviate from the normal representations learned in the training process [5, 6, 112]. In ADBench, semi-supervision refers to *incomplete label learning* in weak-supervision (see [123]). More discussions on semi-supervised AD are deferred to Appx. §A.3.

Table 1: Comparison among ADBench and existing benchmarks, where ADBench comprehensively includes the most datasets and algorithms, uses both real-world and synthetic datasets, covers both shallow and deep learning (DL) algorithms, and consider multiple comparison angles.

Benchmark	Coverage (§3.2)		Data Source		Algorithm Type		Comparison Angle (§3.3)		
	# datasets	# algo.	Real-world	Synthetic	Shallow	DL	Supervision	Types	Robustness
Ruff et al. [90]	3	9	✓	✓	✓	✓	✗	✓	✗
Goldstein et al. [30]	10	19	✓	✗	✓	✗	✗	✓	✗
Domingues et al. [23]	15	14	✓	✗	✓	✗	✗	✗	✓
Soenen et al. [98]	16	6	✓	✗	✓	✗	✗	✗	✗
Steinbuss et al. [99]	19	4	✗	✓	✓	✗	✗	✓	✗
Emmott et al. [25]	19	8	✓	✓	✓	✗	✗	✓	✓
Campos et al. [14]	23	12	✓	✗	✓	✗	✗	✗	✗
ADBench (ours)	55	30	✓	✓	✓	✓	✓	✓	✓

2.2 Existing Datasets and Benchmarks for Tabular AD

AD Datasets in Literature. Existing benchmarks mainly evaluate on a part of the datasets derived from the ODDS Library [86], DAMI Repository [14], and Anomaly Detection Meta-Analysis Benchmarks [25]. In ADBench, we include almost all publicly available datasets, and (for the first time) add larger datasets adapted from CV and NLP domains, for a more holistic view. See §3.2.

Existing Benchmarks. There are some notable works that take effort to benchmark AD methods on tabular data, e.g., [14, 23, 25, 90, 99] (see Appx. A.4). How does ADBench differ from them?

First, previous studies mainly focus on benchmarking the shallow unsupervised AD methods. Considering the rapid advancement of ensemble learning and deep learning methods, we argue a comprehensive benchmark should consider them as well. Second, most existing works only evaluate on public real-world datasets and/or some fully synthetic datasets, we organically incorporate both of them to unlock deeper insights. More importantly, existing benchmarks primarily focus on direct performance comparisons, while the settings may not be sufficiently complex to understand AD algorithm characteristics. We strive to address the above issues in ADBench, and illustrate the main differences among the proposed ADBench and existing AD benchmarks in Table 1.

Note that the term “anomaly detection” is overloaded in many fields; there are some AD benchmarks for time-series [53], CV [15, 121] and NLP [84], but they are different from tabular AD in nature.

3 ADBench: AD Benchmark Driven by Research and Application Needs

3.1 Preliminaries and Problem Definition

Unsupervised AD often presents a collection of n samples $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathbb{R}^{n \times d}$, where each sample has d features. Given the inductive setting, the goal is to train an AD model M to output anomaly score $\mathbf{O} := M(\mathbf{X}) \in \mathbb{R}^{n \times 1}$, where higher scores denote for more outlyingness. In the inductive setting, we need to predict on $\mathbf{X}_{\text{test}} \in \mathbb{R}^{m \times d}$, so to return $\mathbf{O}_{\text{test}} := M(\mathbf{X}_{\text{test}}) \in \mathbb{R}^{m \times 1}$.

Supervised AD also has the (binary) ground truth labels of \mathbf{X} , i.e., $\mathbf{y} \in \mathbb{R}^{n \times 1}$. A supervised AD model M is first trained on $\{\mathbf{X}, \mathbf{y}\}$, and then returns anomaly scores for the $\mathbf{O}_{\text{test}} := M(\mathbf{X}_{\text{test}})$.

Semi-supervised AD only has the partial label information $\mathbf{y}^l \in \mathbf{y}$. The AD model M is trained on the entire feature space \mathbf{X} with the partial label \mathbf{y}^l , i.e., $\{\mathbf{X}, \mathbf{y}^l\}$, and then outputs $\mathbf{O}_{\text{test}} := M(\mathbf{X}_{\text{test}})$.

Remark. Irrespective of the types of underlying AD algorithms, the goal of ADBench is to understand AD algorithms’ performance under the inductive setting. Collectively, we refer semi-supervised and supervised AD methods as “label-informed” methods. Refer to §4.1 for specific experiment settings.

3.2 The Largest AD Benchmark with 30 Algorithms and 55 Datasets

Algorithms. Compared to the previous studies, we have a larger algorithm collection with (i) latest unsupervised AD algorithms like DeepSVDD [91] and ECOD [60]; (ii) SOTA semi-supervised algorithms, including DeepSAD [92] and DevNet [78]; (iii) latest network architectures like ResNet [38] in computer vision (CV) and Transformer [103] in natural language processing (NLP) domain—we adapt ResNet and FTTransformer models [33] for tabular AD in the proposed ADBench; and (iv) ensemble learning methods like LightGBM [46], XGBoost [17], and CatBoost [81]. Fig. 2 shows the algorithms (14 unsupervised, 7 semi-supervised, and 9 supervised algorithms) in ADBench.

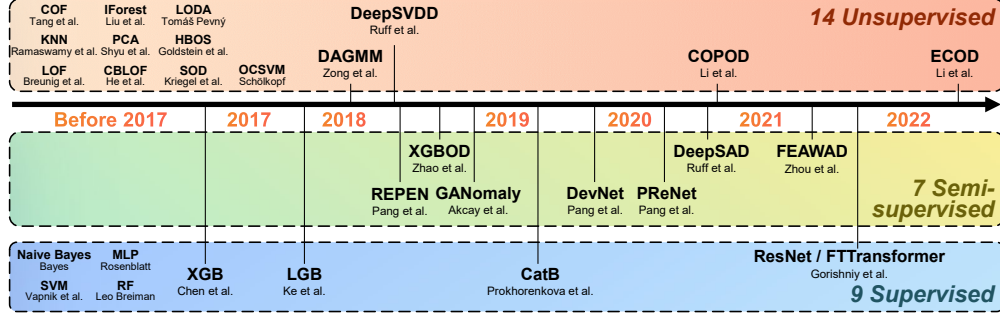


Figure 2: ADBench covers a wide range of AD algorithms. See Appx. B.2 for more details.

Algorithm Implementation. Most unsupervised algorithms are readily available in our early work Python Outlier Detection (PyOD) [118], and some supervised methods are available in scikit-learn [79] and corresponding libraries. Supervised ResNet and FTTransformer tailored for tabular data have been open-sourced in their original paper [33]. We implement all the semi-supervised methods and release them along with ADBench—we consider them as great addition to the community.

Public AD Datasets. In ADBench, we gather more than 40 public real-world datasets for model evaluation, as shown in Appx. Table 2. These datasets cover many application domains, including healthcare (e.g., disease diagnosis), audio and language processing (e.g., speech recognition), image processing (e.g., object identification), finance (e.g., financial fraud detection), etc.

Newly-added Datasets in ADBench. As most of these datasets are relatively small, we introduce 7 more complex datasets from CV and NLP domains with more samples and richer features in ADBench (highlighted in Appx. Table 2). Pretrained models are applied to extract data embedding from NLP and CV datasets to access more complex representation. For NLP datasets, we use BERT [47] pretrained on the BookCorpus and English Wikipedia to extract the embedding of the [CLS] token. For CV datasets, we use ResNet18 [38] pretrained on the ImageNet [22] to extract the embedding after the last average pooling layer. Following previous works [91, 92], we set one of the multi-classes as normal and downsample the remaining classes to 5% of the total instances as anomalies, and report the average results over all the respective classes.

3.3 Benchmark Angles in ADBench

3.3.1 Angle I: Availability of Ground Truth Labels (Supervision)

Motivation. As shown in Table 1, existing benchmarks only focus on the unsupervised setting, i.e., none of the labeled anomalies is available. Despite, in addition to numerous unlabeled samples, one may have access to a limited number of labeled anomalies in real-world applications, e.g., a few anomalies identified by domain experts or human-in-the-loop techniques like active learning, which shows their importance in many works [4, 5, 50, 113]. Notably, there are a group of semi-supervised AD algorithms [74, 75, 77, 78, 92, 101, 122], that have not been covered by existing benchmarks.

Our design: We first benchmark existing unsupervised anomaly detection methods, and then evaluate both semi-supervised and fully-supervised methods with varying level of supervision following the settings in [74, 78, 122] to provide a fair comparison. For example, labeled anomalies $\gamma_l = 10\%$ means that 10% anomalies in the training set is known while other samples remain unlabeled. The full experiment results of un-, semi-, and full-supervised algorithms are presented in §4.2.

3.3.2 Angle II: Types of Anomalies

Motivation. While extensive real-world datasets can be used for benchmarking, they often consist of a mixture of different types of anomalies, making it challenging to understand the pros and cons of AD algorithms regarding specific anomaly types [32, 99]. In real-world applications, one may know specific types of anomalies at interest in prior. To better understand the impact of anomaly types, we create synthetic datasets based on real-world datasets by injecting specific types of anomalies to analyze the response of AD algorithms.

Our design: In ADBench, we create *realistic* synthetic datasets from real-world datasets by injecting specific types of anomalies. Some existing works, such as PyOD [118], simply generate fully

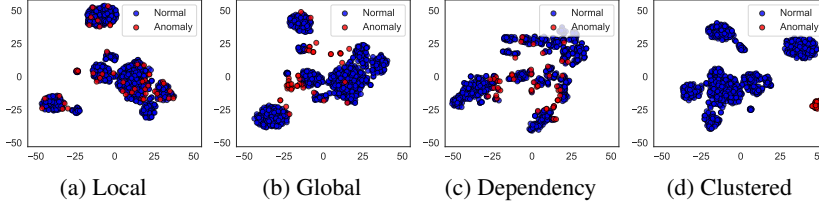


Figure 3: Illustration of four types of synthetic anomalies shown on Lymphography dataset.

synthetic anomalies by assuming their data distribution, which fail to generate complex anomalies. We follow and enrich the approach in [99] to generate “realistic” synthetic data; ours supports more types of anomaly generation. The core idea is to build a generative model (e.g., Gaussian mixture model GMM used in both [99] and ADBench) using the normal samples from a real-world dataset and discard all its original anomalies as we do not know their types. With the generative model, we could then generate normal samples, and also generate different types of anomalies based on their definitions by tweaking the generative model. Thus, the generation of normal samples are the same in all settings if not noted, and we provide the details of various types of anomalies as below.

Definition and Generation Process of Four Types of Common Anomalies Used in ADBench:

- **Local anomalies** refer to the anomalies that are deviant from their local neighborhoods [12]. We follow the GMM procedure [69, 99] 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 more different from the normal data [41], generated from a uniform distribution $\text{Unif}(\alpha \cdot \min(\mathbf{X}^k), \alpha \cdot \max(\mathbf{X}^k))$, where the boundaries are defined as the *min* and *max* of an input feature, e.g., k -th feature \mathbf{X}^k , and $\alpha = 1.1$ controls the outlyingness of anomalies.
- **Dependency anomalies** refer to the samples that do not follow the dependency structure which normal data follow [68], i.e., the input features of dependency anomalies are assumed to be independent to each other. Vine Copula [1] method is applied to model the dependency structure of original data, where the probability density function of generated anomalies are set to complete independence by removing the modeled dependency (see [68]). We use Kernel Density Estimation (KDE) [37] to estimate the probability density function of features and generate normal samples.
- **Clustered anomalies**, also known as group anomalies [56], exhibit similar characteristics [25, 62]. We scale the mean feature vector of normal samples by $\alpha = 5$, i.e., $\hat{\mu} = \alpha \hat{\mu}$, where α controls the distance between anomaly clusters and the normal, and use the scaled GMM to generate anomalies.

Fig. 3 shows 2-d t-SNE [102] visualization of the four types of synthetic outliers generated from Lymphography dataset, where they generally satisfy the expected characteristics. Local anomalies (Fig. 3a) are well overlapped with the normal samples. Global anomalies (Fig. 3b) are more deviated from the normal samples and on the edges of normal clusters. The other two types of anomalies are as expected with no clear dependency structure in Fig. 3c and having anomaly cluster(s) in Fig. 3d. In ADBench, we analyze the algorithm performances under all four types of anomalies above (§4.3).

3.3.3 Angle III: Model Robustness with Noisy and Corrupted Data

Motivation. Model robustness has been an important aspect of anomaly detection and adversarial machine learning [13, 24, 26, 48, 107]. Meanwhile, it is likely that the input data suffers from noise and corruption to some extent in real-world applications [25, 32, 36, 71]. However, this important view has not been well studied in existing benchmarks, and we try to understand this by evaluating AD algorithms under three noisy and corruption settings (see results in §4.4):

- **Duplicated Anomalies.** In real-world applications, it is likely that certain anomalies repeat multiple times in the data, and the presence of duplicated anomalies is also called the “anomaly masking” [32, 36, 63], posing challenges to many algorithms [14], e.g., the density-based KNN [7, 85]. Besides, the change of anomaly frequency would also affect the behavior of detection methods [25]. Therefore, we simulate this noise by duplicating the anomalies (both features and labels) up to 6 times, and observe how do AD algorithms respond to it.
- **Irrelevant Features.** Tabular data may contain irrelevant features caused by measurement noise or inconsistent measuring units [16, 32], where these noisy dimensions could hide the characteristics of anomaly data and thus make the detection process more difficult [75, 90]. We add irrelevant features up to 50% of the total input features (i.e., d in the problem definition) by generating

uniform noise features from $\text{Unif}(\min(\mathbf{X}^k), \max(\mathbf{X}^k))$ of randomly selected k -th input feature \mathbf{X}^k while the labels stay correct, and summarize the algorithm performance changes.

- **Annotation Errors.** While existing studies [78, 92] explored anomaly contamination in the unlabeled samples, we further discuss more generalized impact of label contamination on the algorithm performance, where the label flips [70, 120] between the normal samples and anomalies are considered (up to 50% of total labels). Note this setting does not affect unsupervised methods as they do not use any labels. Discussion of annotation errors is meaningful since manual annotation or some automatic labeling techniques are always noisy while being treated as perfect.

4 Experiment Results and Analyses

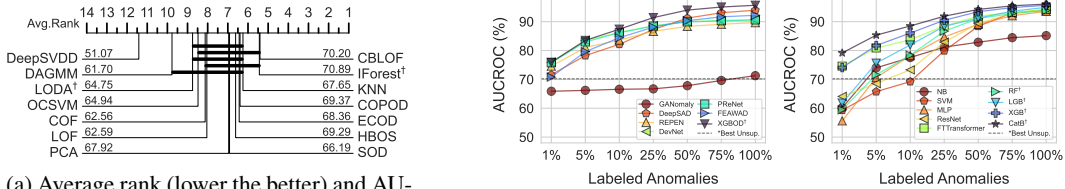
We conduct 93,654 experiments (Appx. C) to answer **Q1** (§4.2): How do AD algorithms perform under real-world datasets with varying levels of supervision? **Q2** (§4.3): How do AD algorithms respond to different types of anomalies? **Q3** (§4.4): How robust are AD algorithms with noisy and corrupted data? In each subsection, we first present the key results and analyses (please refer to the additional points in Appx. D), and then propose a few open questions and future research directions.

4.1 Experiment Setting

Datasets, Train/test Data Split, and Independent Trials. As described in §3.2 and Appx. Table 2, ADBench includes 55 existing and freshly proposed datasets, which cover different fields including healthcare, security, and more. Although unsupervised AD algorithms are primarily designed for the transductive setting (i.e., outputting the anomaly scores on the input data only other than making predictions on newcoming data), we adapt all the algorithms for the inductive setting to make prediction on the newcoming data, which is helpful in applications and also common in popular AD library PyOD [118], TODS [52], and PyGOD [64]. Thus, we use 70% data for training and the remaining 30% as testing set. We use stratified sampling to keep the anomaly ratio consistent. We repeat each experiment 3 times and report the average. Detailed settings are described in Appx. C.

Hyperparameter Settings. For all the algorithms in ADBench, we use their default hyperparameter (HP) settings in the original paper for fair comparison. Refer to the Appx. C for more information.

Evaluation Metrics and Statistical Tests. We evaluate different AD methods by two widely used metrics: AUCROC (Area Under Receiver Operating Characteristic Curve) and AUCPR (Area Under Precision-Recall Curve) value¹. Besides, the critical difference diagram (CD diagram) [21, 43] based on the Wilcoxon-Holm method is used for comparing groups of AD methods statistically ($p \leq 0.05$).



(a) Average rank (lower the better) and AUCROC of unsupervised methods with CD diagram; groups of algorithms not significantly different are connected horizontally.

(b) AUCROC vs. % of labeled anomalies (x-axis); semi-supervised (left) and fully-supervised (right). The best unsupervised algorithm CBLOF is denoted as the dashed line. "†" marks ensembling.

Figure 4: AD model performance on 55 real-world datasets. (a) shows that no unsupervised algorithm can statistically outperform. (b) shows that semi-supervised methods leverage the labels more efficiently than fully-supervised methods with small labeled anomaly ratio γ_l .

4.2 Overall Model Performance on Real-world Datasets with Varying Degrees of Supervision

As layed out in §3.3.1, we first present the results of unsupervised methods on 55 datasets in Fig. 4a, and then compare label-informed semi- and fully-supervised methods under varying degrees of supervision, i.e., different label ratios of γ_l (from 1% to 100% full labeled anomalies) in Fig. 4b.

None of the unsupervised methods is statistically better than the others, as shown in the CD diagram of Fig. 4a. We also note that DL-based unsupervised methods like DeepSVDD and DAGMM

¹We present the results based on AUCROC and observe similar results for AUCPR; See Appx. D for all.

are surprisingly worse than shallow methods. Without the guidance of label information, DL-based unsupervised algorithms are often harder to train (due to more hyperparameters), as well as more difficult to tune their hyperparameters, leading to unsatisfactory performance.

Semi-supervised methods outperform supervised methods when limited label information is available. For $\gamma_l \leq 5\%$, i.e., only less than 5% labeled anomalies are available during training, the detection performance of semi-supervised methods (median AUCROC= 74.59% for $\gamma_l = 1\%$ and AUCROC= 81.04% for $\gamma_l = 5\%$) are generally better than that of fully-supervised algorithms (median AUCROC= 61.80% for $\gamma_l = 1\%$ and AUCROC= 74.00% for $\gamma_l = 5\%$). For most of the semi-supervised methods, merely 1% labeled anomalies are sufficient to surpass the best unsupervised method (shown as the dashed line in Fig. 4b), while most of supervised methods need 10% labeled anomalies to achieve so. We also show the improvement of algorithm performances with regard to the increasing γ_l , and we notice that with large amount of labeled anomalies, both semi-supervised and supervised methods have close performance.

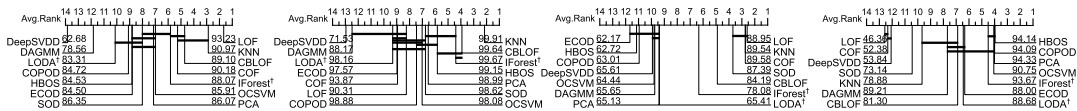
Latest network architectures like Transformer and emerging ensemble methods yield competitive performance in AD. Fig. 4b shows FTTransformer and ensemble methods like XGB(oost) and CatB(oost) provide satisfying detection performance among all the label-informed algorithms, even these methods are not specifically proposed for the anomaly detection tasks. For $\gamma_l = 1\%$, the AUCROC of FTTransformer and the median AUCROC of ensemble methods are 74.61% and 74.08%, respectively, outperform the median AUCROC of all label-informed methods 71.19%.

Future Direction 1: Unsupervised Algorithm Evaluation, Selection, and Design. For unsupervised AD, the results suggest that future algorithms should be evaluated on large testbeds like ADBench for statistical tests (such as via critical different diagram). Meanwhile, no-free-lunch theorem [106] suggests there is no universal winner for all tasks, and more focus should be spent on understanding the suitability of each AD algorithm. Specifically, algorithm selection is important in unsupervised AD, but limited works [8, 67, 119] have studied this. We may consider self-supervision [83, 95, 108] and transfer learning [20] to improve tabular AD as well. Thus, we suggest more focus on large-scale evaluation, task-driven algorithm selection, and data augmentation/transfer for unsupervised AD.

Future Direction 2: Semi-supervised Learning. By observing the success of using limited labels in AD, we would call for more attention on semi-supervised AD methods which can leverage both the guidance from labels efficiently and the exploration from the unlabeled data. Regarding to backbones, latest network architectures like Transformer and ensembling show their superiority in AD tasks.

4.3 Algorithm Performance under Different Types of Anomalies

Under four types of anomalies introduced in §3.3.2), we show the performances of unsupervised methods in Fig. 5, and then compare both semi- and fully-supervised methods in Fig. 6.



(a) Local anomalies (b) Global anomalies (c) Dependency anomalies (d) Clustered anomalies
Figure 5: Avg. rank (lower the better) of unsupervised methods on different types of anomalies. Groups of algorithms not significantly different are connected horizontally in the CD diagrams. The unsupervised methods perform well when their assumptions conform to the anomaly types.

Performance of unsupervised algorithms highly depends on the alignment of its assumptions and the underlying anomaly type. As expected, *local* anomaly factor (LOF) is statistically better than other unsupervised methods for the local anomalies (Fig. 5a), and KNN, which uses *k*-th (*global*) nearest neighbor’s distance as anomaly scores, is the statistically best detector for global anomalies (Fig. 5b). Again, there is no algorithm performing well on all types of anomalies; KNN achieves the best AUCROC on global anomalies (Fig. 5b) and the second best AUCROC on dependency anomalies (Fig. 5c), but performs poorly on the clustered anomaly (Fig. 5d). Practitioners should select algorithms based on the characteristics of underlying task, and consider the algorithm which may cover more high-interest anomaly types [56].

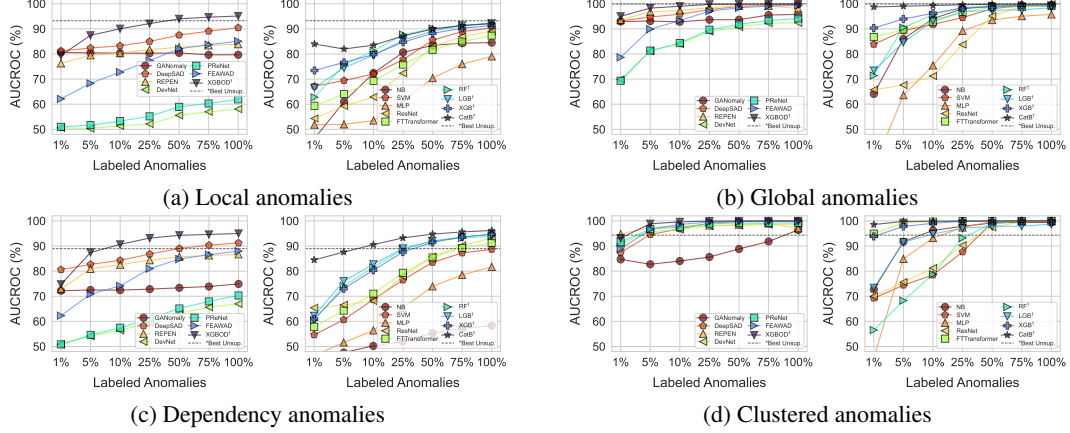


Figure 6: Semi- (left of each subfigure) and supervised (right) algorithms’ performance on different types of anomalies with varying level of labeled anomalies. Surprisingly, these label-informed algorithms are *inferior* to the best unsupervised method except for the clustered anomalies.

The “power” of prior knowledge on anomaly types may outweigh the usage of partial labels.

For the local, global, and dependency anomalies, where most of the label-informed methods perform worse than the best unsupervised methods of each type (corresponding to LOF, KNN, and LOF). For example, the detection performance of XGBOD for the local anomalies is inferior to the best unsupervised method LOF when $\gamma_l \leq 50\%$, while other methods perform worse than LOF in all cases (See Fig. 6a). Why could not label-informed algorithms beat unsupervised methods in this setting? We believe that partially labeled anomalies cannot well capture all characteristics of specific types of anomalies, and learning such decision boundaries is challenging. For instance, different local anomalies often exhibit various behaviors, as shown in Fig. 3, which may be easier to identify by a generic definition of “locality” in unsupervised methods other than specific labels. Thus, incomplete label information may bias the learning process of these label-informed methods, which explains their relatively inferior performances compared to best unsupervised methods. This conclusion is further verified by the results of clustered anomalies (See Fig. 6d), where label-informed (especially semi-supervised) methods outperform the best unsupervised method HBOS, as few labeled anomalies can already represent the similar behaviors in the clustered anomalies (Fig. 3d).

Future Direction 3: Leveraging Anomaly Types as Valuable Prior Knowledge. The above results emphasize the importance of knowing anomaly types in achieving high detection performance even without using any labels, and calls for attention to design anomaly-type-aware detection algorithms. In an ideal world, one may combine multiple AD algorithms based on the composition of anomaly types, via frameworks like dynamic model selection and combination [117]. To our knowledge, the latest advancement in this end [44] provides an equivalence criterion for measuring to what degree two anomaly detection algorithms detect the same kind of anomalies. Furthermore, future research may also consider designing semi-supervised AD methods capable of detecting different types of unknown anomalies while effectively improving performance by the partially available labeled data.

4.4 Algorithm Robustness under Noisy and Corrupted Data

In this section, we investigate the algorithm robustness (i.e., Δ performance) of different AD algorithms under noisy and data corruption described in §3.3.3. The default γ_l is set to 100% since we only care about the relative change of model performance. Fig. 7 demonstrates the results.

Unsupervised methods are more susceptible to duplicated anomalies. As shown in Fig. 7a, almost all unsupervised methods are severely impacted by duplicated anomalies. Their AUCROC deteriorates proportionally with the increase of duplication. When anomalies are duplicated by 6 times, the median Δ AUCROC of unsupervised methods is -17.42% , compared to that of semi-supervised methods -0.04% (Fig. 7b) and supervised methods 0.54% (Fig. 7c). One explanation is that unsupervised methods often assume the underlying data is imbalanced with only a smaller percentage of anomalies—they rely on this assumption to detect anomalies. With more duplicated anomalies, the underlying data becomes more balanced and the minority assumption of anomalies is

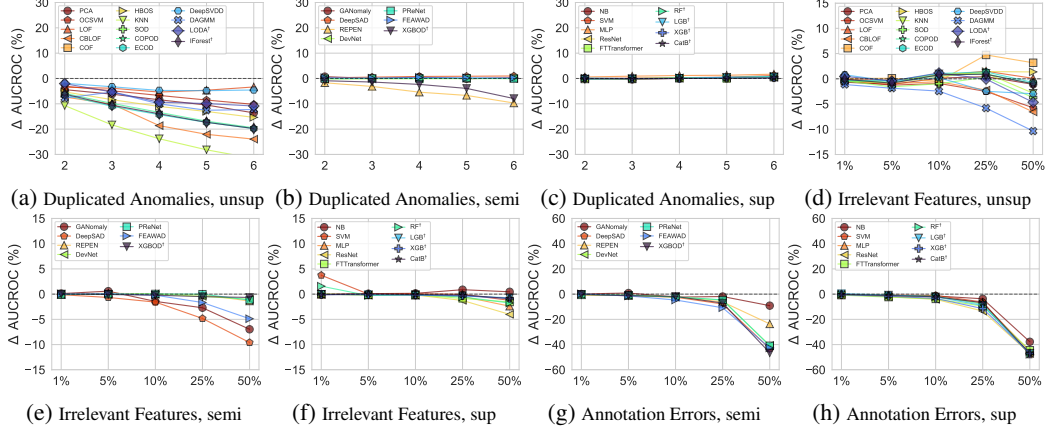


Figure 7: Algorithm performance change under noisy and corrupted data (i.e., duplicated anomalies for (a)-(c), irrelevant features for (d)-(f), and annotation errors for (g) and (h)). x-axis denotes either the duplicated times or noise ratio. y-axis denotes the % of performance change (Δ AUCROC) and its range remains consistent across different algorithms. The results reveals unsupervised methods’ susceptibility to duplicated anomalies and the usage of label information in defending irrelevant features. Un-, semi-, and fully-supervised methods are denoted as *unsup*, *semi*, and *sup*, respectively.

violated, causing the degradation of unsupervised methods. Differently, more balanced datasets do not affect the performance of semi- and fully-supervised methods remarkably, with the help of labels.

Irrelevant features cause little impact on supervised methods due to feature selection. Compared to the unsupervised and most semi-supervised methods, the training process of supervised methods is fully guided by the data labels (y), therefore perform robustly to the irrelevant features (i.e., corrupted X) due to the direct (or indirect) feature selection process. For instance, ensemble trees like XGBoost can filter irrelevant features. As shown in Fig. 7f, even the worst supervised algorithm (say ResNet) in this setting yields $\leq 5\%$ degradation when 50% of the input features are corrupted by the uniform noises, while the un- and semi-supervised methods could face up to 10% degradation. Besides, the robust performances of supervised methods (and some semi-supervised methods like DevNet) indicate that the label information can be beneficial for feature selection. Also, Fig. 7f shows minor irrelevant features (e.g., 1%) even help supervised methods as regularization to generalize better.

Both semi- and fully-supervised methods shows great resilience to minor annotation errors. Although the detection performance of these methods significantly downgrade when the annotation errors are severe (as shown in Fig. 7g and 7h), their degradation with regard to minor annotation errors is acceptable. The median Δ AUCROC of semi- and fully-supervised methods for 5% annotation errors is -2.25% and -2.2% , respectively. That being said, label-informed methods are still acceptable in practice as the annotation error should be relatively small [58, 109].

Future Direction 4: Noise-resilient AD Algorithms. Our results indicate there is an improvement space for robust unsupervised AD algorithms. One immediate remedy is to incorporate unsupervised feature selection [18, 72, 73] to combat irrelevant features. Moreover, label information could serve as an effective guidance of model training against data noise, and it helps semi- and fully-supervised methods to be more robust. Given the difficulty of acquiring full labels, we would suggest considering semi-supervised methods as the backbone for designing more robust AD algorithms.

5 Conclusions and Future Work

In this paper, we introduce ADBench, the most comprehensive tabular anomaly detection benchmark with 30 algorithms and 55 benchmark datasets. Based on the analyses on multiple comparison angles, we unlock insights on the role of supervision, the importance of prior knowledge on anomaly types, and the principles of designing robust detection algorithms. On top of them, we summarize a few promising future research directions for anomaly detection, along with the fully released benchmark suite for evaluation on new algorithms. ADBench can extend to understand the algorithm response under mixed types of anomalies, include hyperparameter tuning with labels, and enrich by datasets from emerging fields like drug discovery [42], molecule optimization [27, 28], and machine bias [40].

References

- [1] K. Aas, C. Czado, A. Frigessi, and H. Bakken. Pair-copula constructions of multiple dependence. *Insurance: Mathematics and economics*, 44(2):182–198, 2009.
- [2] C. C. Aggarwal. An introduction to outlier analysis. In *Outlier analysis*, pages 1–34. Springer, 2017.
- [3] C. C. Aggarwal. *Neural Networks and Deep Learning - A Textbook*. Springer, 2018.
- [4] N. B. Aissa and M. Guerroumi. Semi-supervised statistical approach for network anomaly detection. *Procedia Computer Science*, 83:1090–1095, 2016.
- [5] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *Asian conference on computer vision*, pages 622–637. Springer, 2018.
- [6] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon. Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019.
- [7] F. Angiulli and C. Pizzuti. Fast outlier detection in high dimensional spaces. In *European conference on principles of data mining and knowledge discovery*, pages 15–27. Springer, 2002.
- [8] M. Bahri, F. Salutari, A. Putina, and M. Sozio. Automl: state of the art with a focus on anomaly detection, challenges, and research directions. *International Journal of Data Science and Analytics*, pages 1–14, 2022.
- [9] T. Bayes. Lii. an essay towards solving a problem in the doctrine of chances. by the late rev. mr. bayes, frs communicated by mr. price, in a letter to john canton, amfr s. *Philosophical transactions of the Royal Society of London*, (53):370–418, 1763.
- [10] L. Bergman and Y. Hoshen. Classification-based anomaly detection for general data. In *International Conference on Learning Representations*, 2019.
- [11] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [12] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. Lof: identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 93–104, 2000.
- [13] H. Cai, J. Liu, and W. Yin. Learned robust pca: A scalable deep unfolding approach for high-dimensional outlier detection. *Advances in Neural Information Processing Systems*, 34, 2021.
- [14] G. O. Campos, A. Zimek, J. Sander, R. J. Campello, B. Micenková, E. Schubert, I. Assent, and M. E. Houle. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data mining and knowledge discovery*, 30(4):891–927, 2016.
- [15] R. Chan, K. Lis, S. Uhlemeyer, H. Blum, S. Honari, R. Siegwart, P. Fua, M. Salzmann, and M. Rottmann. Segmentmeifyoucan: A benchmark for anomaly segmentation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [16] C.-H. Chang, J. Yoon, S. Arik, M. Udell, and T. Pfister. Data-efficient and interpretable tabular anomaly detection. *arXiv preprint arXiv:2203.02034*, 2022.
- [17] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [18] L. Cheng, Y. Wang, X. Liu, and B. Li. Outlier detection ensemble with embedded feature selection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3503–3512, 2020.
- [19] C. Cortes and V. Vapnik. Support vector machine. *Machine learning*, 20(3):273–297, 1995.
- [20] L. Deecke, L. Ruff, R. A. Vandermeulen, and H. Bilen. Transfer-based semantic anomaly detection. In *International Conference on Machine Learning*, pages 2546–2558. PMLR, 2021.
- [21] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research*, 7:1–30, 2006.
- [22] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

- [23] R. Domingues, M. Filippone, P. Michiardi, and J. Zouaoui. A comparative evaluation of outlier detection algorithms: Experiments and analyses. *Pattern Recognition*, 74:406–421, 2018.
- [24] X. Du, J. Zhang, B. Han, T. Liu, Y. Rong, G. Niu, J. Huang, and M. Sugiyama. Learning diverse-structured networks for adversarial robustness. In *International Conference on Machine Learning*, pages 2880–2891. PMLR, 2021.
- [25] A. Emmott, S. Das, T. Dietterich, A. Fern, and W.-K. Wong. A meta-analysis of the anomaly detection problem. *arXiv preprint arXiv:1503.01158*, 2015.
- [26] L. Fan, S. Liu, P.-Y. Chen, G. Zhang, and C. Gan. When does contrastive learning preserve adversarial robustness from pretraining to finetuning? *Advances in Neural Information Processing Systems*, 34, 2021.
- [27] T. Fu, C. Xiao, X. Li, L. M. Glass, and J. Sun. Mimosa: Multi-constraint molecule sampling for molecule optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 125–133, 2021.
- [28] T. Fu, C. Xiao, and J. Sun. Core: Automatic molecule optimization using copy & refine strategy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 638–645, 2020.
- [29] M. Goldstein and A. Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: poster and demo track*, 9, 2012.
- [30] M. Goldstein and S. Uchida. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PloS one*, 11(4):e0152173, 2016.
- [31] I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016.
- [32] P. Gopalan, V. Sharan, and U. Wieder. Pidforest: anomaly detection via partial identification. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] Y. Gorishniy, I. Rubachev, V. Khurlov, and A. Babenko. Revisiting deep learning models for tabular data. *Advances in Neural Information Processing Systems*, 34, 2021.
- [34] N. Görnitz, M. Kloft, K. Rieck, and U. Brefeld. Toward supervised anomaly detection. *Journal of Artificial Intelligence Research*, 46:235–262, 2013.
- [35] C. Grunau and V. Rozhoň. Adapting k-means algorithms for outliers. *arXiv preprint arXiv:2007.01118*, 2020.
- [36] S. Guha, N. Mishra, G. Roy, and O. Schrijvers. Robust random cut forest based anomaly detection on streams. In *International conference on machine learning*, pages 2712–2721. PMLR, 2016.
- [37] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- [38] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [39] Z. He, X. Xu, and S. Deng. Discovering cluster-based local outliers. *Pattern recognition letters*, 24(9-10):1641–1650, 2003.
- [40] X. Hu, Y. Huang, B. Li, and T. Lu. Uncovering the source of machine bias. *27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21), Machine Learning for Consumers and Markets Workshop*, 2021.
- [41] H. Huang, H. Qin, S. Yoo, and D. Yu. Physics-based anomaly detection defined on manifold space. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 9(2):1–39, 2014.
- [42] K. Huang, T. Fu, W. Gao, Y. Zhao, Y. Roohani, J. Leskovec, C. Coley, C. Xiao, J. Sun, and M. Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. *Advances in neural information processing systems*, 2021.
- [43] H. Ismail Fawaz, G. Forestier, J. Weber, L. Idoumghar, and P.-A. Muller. Deep learning for time series classification: a review. *Data mining and knowledge discovery*, 33(4):917–963, 2019.
- [44] C. I. Jerez, J. Zhang, and M. R. Silva. On equivalence of anomaly detection algorithms. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 2022.

- [45] Y. Kawachi, Y. Koizumi, and N. Harada. Complementary set variational autoencoder for supervised anomaly detection. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2366–2370. IEEE, 2018.
- [46] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30:3146–3154, 2017.
- [47] J. D. M.-W. C. Kenton and L. K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186, 2019.
- [48] M. Kim, J. Tack, and S. J. Hwang. Adversarial self-supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:2983–2994, 2020.
- [49] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [50] B. R. Kiran, D. M. Thomas, and R. Parakkal. An overview of deep learning based methods for unsupervised and semi-supervised anomaly detection in videos. *Journal of Imaging*, 4(2):36, 2018.
- [51] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. Outlier detection in axis-parallel subspaces of high dimensional data. In *Pacific-asia conference on knowledge discovery and data mining*, pages 831–838. Springer, 2009.
- [52] K.-H. Lai, D. Zha, G. Wang, J. Xu, Y. Zhao, D. Kumar, Y. Chen, P. Zumkhawaka, M. Wan, D. Martinez, et al. Tods: An automated time series outlier detection system. In *Proceedings of the aaai conference on artificial intelligence*, volume 35, pages 16060–16062, 2021.
- [53] K.-H. Lai, D. Zha, J. Xu, Y. Zhao, G. Wang, and X. Hu. Revisiting time series outlier detection: Definitions and benchmarks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- [54] A. Lazarevic, L. Ertoz, V. Kumar, A. Ozgur, and J. Srivastava. A comparative study of anomaly detection schemes in network intrusion detection. In *Proceedings of the 2003 SIAM international conference on data mining*, pages 25–36. SIAM, 2003.
- [55] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [56] M.-C. Lee, S. Shekhar, C. Faloutsos, T. N. Hutson, and L. Iasemidis. Gen 2 out: Detecting and ranking generalized anomalies. In *2021 IEEE International Conference on Big Data (Big Data)*, pages 801–811. IEEE, 2021.
- [57] M.-C. Lee, Y. Zhao, A. Wang, P. J. Liang, L. Akoglu, V. S. Tseng, and C. Faloutsos. Autoaudit: Mining accounting and time-evolving graphs. In *2020 IEEE International Conference on Big Data (Big Data)*, pages 950–956. IEEE, 2020.
- [58] G. Li, Y. Xie, and L. Lin. Weakly supervised salient object detection using image labels. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [59] Z. Li, Y. Zhao, N. Botta, C. Ionescu, and X. Hu. Copod: copula-based outlier detection. In *2020 IEEE International Conference on Data Mining (ICDM)*, pages 1118–1123. IEEE, 2020.
- [60] Z. Li, Y. Zhao, X. Hu, N. Botta, C. Ionescu, and G. Chen. Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2022.
- [61] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [62] B. Liu, P.-N. Tan, and J. Zhou. Unsupervised anomaly detection by robust density estimation. 2022.
- [63] F. T. Liu, K. M. Ting, and Z.-H. Zhou. Isolation forest. In *2008 eighth IEEE international conference on data mining*, pages 413–422. IEEE, 2008.
- [64] K. Liu, Y. Dou, Y. Zhao, X. Ding, X. Hu, R. Zhang, K. Ding, C. Chen, H. Peng, K. Shu, et al. Pygod: A python library for graph outlier detection. *arXiv preprint arXiv:2204.12095*, 2022.
- [65] S. Liu and M. Hauskrecht. Event outlier detection in continuous time. In *International Conference on Machine Learning*, pages 6793–6803. PMLR, 2021.

- [66] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [67] M. Q. Ma, Y. Zhao, X. Zhang, and L. Akoglu. A large-scale study on unsupervised outlier model selection: Do internal strategies suffice? *arXiv preprint arXiv:2104.01422*, 2021.
- [68] R. Martinez-Guerra and J. L. Mata-Machuca. *Fault detection and diagnosis in nonlinear systems*. Springer, 2016.
- [69] G. W. Milligan. An algorithm for generating artificial test clusters. *Psychometrika*, 50(1):123–127, 1985.
- [70] D. T. Nguyen, C. K. Mummadi, T. P. N. Ngo, T. H. P. Nguyen, L. Beggel, and T. Brox. Self: Learning to filter noisy labels with self-ensembling. In *International Conference on Learning Representations*, 2019.
- [71] G. Pang, L. Cao, and L. Chen. Homophily outlier detection in non-iid categorical data. *Data Mining and Knowledge Discovery*, 35(4):1163–1224, 2021.
- [72] G. Pang, L. Cao, L. Chen, and H. Liu. Unsupervised feature selection for outlier detection by modelling hierarchical value-feature couplings. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 410–419. IEEE, 2016.
- [73] G. Pang, L. Cao, L. Chen, and H. Liu. Learning homophily couplings from non-iid data for joint feature selection and noise-resilient outlier detection. In *IJCAI*, pages 2585–2591, 2017.
- [74] G. Pang, L. Cao, L. Chen, and H. Liu. 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*, pages 2041–2050, 2018.
- [75] G. Pang, C. Ding, C. Shen, and A. v. d. Hengel. Explainable deep few-shot anomaly detection with deviation networks. *arXiv preprint arXiv:2108.00462*, 2021.
- [76] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel. Deep learning for anomaly detection: A review. *ACM Computing Surveys (CSUR)*, 54(2):1–38, 2021.
- [77] G. Pang, C. Shen, H. Jin, and A. v. d. Hengel. Deep weakly-supervised anomaly detection. *arXiv preprint arXiv:1910.13601*, 2019.
- [78] G. Pang, C. Shen, and A. van den Hengel. Deep anomaly detection with deviation networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 353–362, 2019.
- [79] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- [80] T. Pevný. Loda: Lightweight on-line detector of anomalies. *Machine Learning*, 102(2):275–304, 2016.
- [81] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, and A. Gulin. Catboost: unbiased boosting with categorical features. *Advances in neural information processing systems*, 31, 2018.
- [82] C. Qiu, A. Li, M. Kloft, M. Rudolph, and S. Mandt. Latent outlier exposure for anomaly detection with contaminated data. *arXiv preprint arXiv:2202.08088*, 2022.
- [83] C. Qiu, T. Pfaffner, M. Kloft, S. Mandt, and M. Rudolph. Neural transformation learning for deep anomaly detection beyond images. In *International Conference on Machine Learning*, pages 8703–8714. PMLR, 2021.
- [84] M. M. Rahman, D. Balakrishnan, D. Murthy, M. Kutlu, and M. Lease. An information retrieval approach to building datasets for hate speech detection. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [85] S. Ramaswamy, R. Rastogi, and K. Shim. Efficient algorithms for mining outliers from large data sets. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 427–438, 2000.
- [86] S. Rayana. ODDS library, 2016.
- [87] Q. Rebjock, B. Kurt, T. Januschowski, and L. Callot. Online false discovery rate control for anomaly detection in time series. *Advances in Neural Information Processing Systems*, 34, 2021.

- [88] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- [89] S. Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [90] L. Ruff, J. R. Kauffmann, R. A. Vandermeulen, G. Montavon, W. Samek, M. Kloft, T. G. Dietterich, and K.-R. Müller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 2021.
- [91] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402, 2018.
- [92] L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K. Müller, and M. Kloft. Deep semi-supervised anomaly detection. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [93] R. Schirrmeister, Y. Zhou, T. Ball, and D. Zhang. Understanding anomaly detection with deep invertible networks through hierarchies of distributions and features. *Advances in Neural Information Processing Systems*, 33:21038–21049, 2020.
- [94] B. Schölkopf, R. C. Williamson, A. J. Smola, J. Shawe-Taylor, J. C. Platt, et al. Support vector method for novelty detection. In *NIPS*, volume 12, pages 582–588. Citeseer, 1999.
- [95] V. Schweg, M. Chiang, and P. Mittal. Ssd: A unified framework for self-supervised outlier detection. In *International Conference on Learning Representations*, 2020.
- [96] L. Shen, Z. Li, and J. Kwok. Timeseries anomaly detection using temporal hierarchical one-class network. *Advances in Neural Information Processing Systems*, 33:13016–13026, 2020.
- [97] M.-L. Shyu, S.-C. Chen, K. Sarinnapakorn, and L. Chang. A novel anomaly detection scheme based on principal component classifier. Technical report, Miami Univ Coral Gables FI Dept of Electrical and Computer Engineering, 2003.
- [98] J. Soenen, E. Van Wolputte, L. Perini, V. Vercruyssen, W. Meert, J. Davis, and H. Blockeel. 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*, pages 1–9. Outlier Detection and Description Organising Committee, 2021.
- [99] G. Steinbuss and K. Böhm. Benchmarking unsupervised outlier detection with realistic synthetic data. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(4):1–20, 2021.
- [100] J. Tang, Z. Chen, A. W.-C. Fu, and D. W. Cheung. Enhancing effectiveness of outlier detections for low density patterns. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 535–548. Springer, 2002.
- [101] B. Tian, Q. Su, and J. Yin. Anomaly detection by leveraging incomplete anomalous knowledge with anomaly-aware bidirectional gans. *arXiv preprint arXiv:2204.13335*, 2022.
- [102] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [103] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [104] Z. Wang, B. Dai, D. Wipf, and J. Zhu. Further analysis of outlier detection with deep generative models. *Advances in Neural Information Processing Systems*, 33:8982–8992, 2020.
- [105] D. H. Wolpert. Stacked generalization. *Neural networks*, 5(2):241–259, 1992.
- [106] D. H. Wolpert and W. G. Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- [107] B. Wu, J. Chen, D. Cai, X. He, and Q. Gu. Do wider neural networks really help adversarial robustness? *Advances in Neural Information Processing Systems*, 34, 2021.
- [108] Z. Xiao, Q. Yan, and Y. Amit. Do we really need to learn representations from in-domain data for outlier detection? *arXiv preprint arXiv:2105.09270*, 2021.
- [109] Y. Xu, J. Ding, L. Zhang, and S. Zhou. Dp-ssl: Towards robust semi-supervised learning with a few labeled samples. *Advances in Neural Information Processing Systems*, 34, 2021.

- [110] W. Yu, J. Li, M. Z. A. Bhuiyan, R. Zhang, and J. Huai. Ring: Real-time emerging anomaly monitoring system over text streams. *IEEE Transactions on Big Data*, 5(4):506–519, 2017.
- [111] M. D. Zeiler. Adadelata: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- [112] H. Zenati, M. Romain, C.-S. Foo, B. Lecouat, and V. Chandrasekhar. Adversarially learned anomaly detection. In *2018 IEEE International Conference on Data Mining (ICDM)*, pages 727–736. IEEE, 2018.
- [113] D. Zha, K.-H. Lai, M. Wan, and X. Hu. Meta-aad: Active anomaly detection with deep reinforcement learning. In *2020 IEEE International Conference on Data Mining (ICDM)*, pages 771–780. IEEE, 2020.
- [114] J. Zhao, X. Liu, Q. Yan, B. Li, M. Shao, and H. Peng. Multi-attributed heterogeneous graph convolutional network for bot detection. *Information Sciences*, 537:380–393, 2020.
- [115] Y. Zhao and M. K. Hryniewicki. Xgbod: improving supervised outlier detection with unsupervised representation learning. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2018.
- [116] Y. Zhao, X. Hu, C. Cheng, C. Wang, C. Wan, W. Wang, J. Yang, H. Bai, Z. Li, C. Xiao, et al. Suod: Accelerating large-scale unsupervised heterogeneous outlier detection. *Proceedings of Machine Learning and Systems*, 3:463–478, 2021.
- [117] Y. Zhao, Z. Nasrullah, M. K. Hryniewicki, and Z. Li. Lscp: Locally selective combination in parallel outlier ensembles. In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pages 585–593. SIAM, 2019.
- [118] Y. Zhao, Z. Nasrullah, and Z. Li. Pyod: A python toolbox for scalable outlier detection. *Journal of Machine Learning Research*, 20:1–7, 2019.
- [119] Y. Zhao, R. Rossi, and L. Akoglu. Automatic unsupervised outlier model selection. *Advances in Neural Information Processing Systems*, 34, 2021.
- [120] G. Zheng, A. H. Awadallah, and S. Dumais. Meta label correction for noisy label learning. *AAAI 2021*, 2021.
- [121] Y. Zheng, X. Wang, Y. Qi, W. Li, and L. Wu. Benchmarking unsupervised anomaly detection and localization. *arXiv preprint arXiv:2205.14852*, 2022.
- [122] Y. Zhou, X. Song, Y. Zhang, F. Liu, C. Zhu, and L. Liu. Feature encoding with autoencoders for weakly supervised anomaly detection. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [123] Z.-H. Zhou. A brief introduction to weakly supervised learning. *National science review*, 5(1):44–53, 2018.
- [124] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International conference on learning representations*, 2018.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] We describe limitations and future works in §5.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A] The better understanding of AD algorithms could facilitate the model deployment, which will lead to positive societal results.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] To facilitate the reproducibility and fast experimental pipeline for the anomaly detection benchmark, we have made all the benchmark datasets and algorithms publicly available with BSD-2 License at <https://github.com/Minqi824/ADBench>, and welcome any customized algorithms to be evaluated via the plug-and-play testbed of ADBench.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] We provide detailed data description of 55 datasets of proposed ADBench in Appx. B.1, and describe the hyperparameter settings of all the 30 algorithms of ADBench in Appx. B.2. Complete experiment settings of proposed ADBench is presented in Appx. C.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Considering the extensive experiments (93.654 in total) involved in ADBench, we mainly demonstrate the average model performances across different datasets, while we also report the critical difference diagram (CD diagram) based on the Wilcoxon-Holm statistical method. Complete experiment settings are presented in §4.1 and Appx. C.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We include the total amount of compute and the computational resources required for ADBench, and further report the runtime comparison of 30 algorithms of ADBench in Appx. C.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes] We release the assets under an inclusive BSD-2 License.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] The new assets of proposed ADBench, including corresponding datasets and code, are available at <https://github.com/Minqi824/ADBench>.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

Supplementary Material for *ADBench: Anomaly Detection Benchmark*

Additional information on related works, datasets, algorithms, and additional experiment settings and results

A Related Works with More Details

In this section, we provide more details on existing AD algorithms and benchmarks, in addition to the primary content discussed in §2.

A.1 Unsupervised Methods

Unsupervised Methods by Assuming Anomaly Data Distributions. Unsupervised AD methods are proposed with different assumptions of data distribution [2], e.g., anomalies locate in low-density regions, and their performance often depends on the agreement between the input data and the algorithm assumption(s). Many unsupervised methods have been proposed in the last few decades [2, 10, 76, 90, 118], which can be roughly categorized into shallow and deep (neural network) methods. The former often carries better interpretability, while the latter handles large, high-dimensional data better. Recent book [2] and surveys [76, 90] provide great details on these algorithms, while we further elaborate on a few representative unsupervised methods. More algorithm details and hyperparameter settings are illustrated in Appx. §B.2

Representative Shallow Methods. Some representative shallow methods include: (i) Isolation Forest (IForest) [63] builds an ensemble of trees to isolate the data points and defines the anomaly score as the distance of an individual instance to the root; (ii) One-class SVM (OCSVM) [94] maximizes the margin between origin and the normal samples, where the decision boundary is the hyper-plane that determines the margin; and (iii) Empirical-Cumulative-distribution-based Outlier Detection (ECOD) [60] computes the empirical cumulative distribution per dimension of the input data, and then aggregates the tail probabilities per dimension for calculating the anomaly score.

Representative Deep Methods. Deep (neural network) methods gain more attention recently, and we briefly review some representative ones in this section. Deep Autoencoding Gaussian Mixture Model (DAGMM) [124] jointly optimizes the deep autoencoder and the Gaussian mixture model in an end-to-end neural network fashion. The joint optimization balances autoencoding reconstruction, density estimation of latent representation, and regularization, and helps the autoencoder escape from less attractive local optima and further reduce reconstruction errors, avoiding the need of pre-training. Deep Support Vector Data Description (DeepSVDD) [91] trains a neural network to learn a transformation that minimizes the volume of a data-enclosing hypersphere in the output space, and calculates the anomaly score as the distance of transformed embedding to the center of the hypersphere.

A.2 Supervised Methods

Due to the difficulty and cost of collecting large-scale labeled data, fully-supervised anomaly detection is often impractical as it assumes the availability of labeled training data with both normal and anomaly samples [76]. Although some loss functions (e.g., the focal loss [61]) are devised to address the class imbalance problem, they are often not specific for AD tasks. There also exist few works [34, 45] discussing the relationship between fully-supervised and semi-supervised AD methods, and argue that semi-supervised AD needs to ground on the unsupervised learning paradigm instead of the supervised one for detecting both known and unknown anomalies. We implement several representative supervised classification algorithms in ADBench (as shown in Appx. §B.2), and recommend interesting readers to recent machine learning books [3, 31] and scikit-learn [79] for more details about recent supervised methods designed for the classification tasks.

A.3 Semi-supervised Methods

Semi-supervised AD algorithms are designed to capitalize the supervision from partial labels, while keeping the ability of detecting unseen types of anomalies. To this end, some recent studies investigate to efficiently use partially labeled data for improving detection performance, and leverage the unlabeled data to facilitate representation learning. We further provide some technical details on representative semi-supervised AD methods here. Please see Appx. §B.2 for more algorithm details and hyperparameter settings in ADBench.

Representative Methods. Extreme Gradient Boosting Outlier Detection (XGBOD) [115] uses multiple unsupervised AD algorithms to extract useful representations from the underlying data that augment the predictive capabilities of an embedded supervised classifier on an improved feature space. Deep Semi-supervised Anomaly Detection (DeepSAD) [92] is an end-to-end methodology considered as the state-of-the-art method in semi-supervised anomaly detection. DeepSAD improves the DeepSVDD [91] model by the inverse loss function for the labeled anomalies. REPresentations for a random nEarest Neighbor distance based metho (REPEN)

[74] proposes a ranking model-based framework, which unifies representation learning and anomaly detection to learn low-dimensional representations tailored for random distance-based detectors. Deviation Networks (DevNet) [78] constructs an end-to-end neural network for learning anomaly score, which forces the network to produce statistically higher anomaly score for identified anomalies than that of unlabeled data. Pairwise Relation prediction-based ordinal regression Network (PReNet) [77] formulates the anomaly detection problem as a pairwise relation prediction task, which defines a two-stream ordinal regression neural network to learn the relation of randomly sampled instance pairs. Feature Encoding with AutoEncoders for Weakly-supervised Anomaly Detection (FEAWAD) [122] leverages an autoencoder to encode the input data and utilize hidden representation, reconstruction residual vector and reconstruction error as the new representations for improving the DevNet [78] and DAGMM [124].

A.4 Existing AD Benchmarks

As we show in Table 1, there are a line of existing AD benchmarks. [90] discusses a unifying review of both the shallow and deep anomaly detection methods, but they mainly focus on the theoretical perspective and thus are lack of results from the experimental views. [14] benchmarks 19 different unsupervised methods on 10 datasets, and analyze the characteristics of density-based and clustering-based algorithms. [23] tests 14 unsupervised anomaly detection methods on 15 real-world datasets, and analyzes the scalability, memory consumption and robustness of different methods. [99] proposes a generic process for the generation of realistic synthetic data. The synthetic normal instances are reconstructed from existing real-world benchmark data, while synthetic anomalies are in line with a characterizable deviation from the modeling of synthetic normal data. [25] evaluates 8 unsupervised methods on 19 real-world datasets, and produce a large corpus of synthetic anomaly detection datasets that vary in their construction across several dimensions that are important to real-world applications. [14] compares 12 unsupervised anomaly detection approaches on 23 datasets, providing characterization of benchmark datasets and their suitability as anomaly detection benchmark sets.

All these existing works lay the foundation of AD algorithm design, and we further improve the foundation by considering more datasets, algorithms, and comparison aspects.

B More Details on ADBench

B.1 ADBench Dataset List

As described in §3.2, ADBench is the largest AD benchmark with 55 datasets. More specifically, Table 2 shows the datasets used in ADBench, covering many application domains, including healthcare (e.g., disease diagnosis), audio and language processing (e.g., speech recognition), image processing (e.g., object identification), finance (e.g., financial fraud detection), and more, where we shown this information in the last column. We resample the sample size to 1,000 for those datasets smaller than 1,000, and use the subsets of 10,000 for those datasets greater than 10,000 for considering the computational cost.

Newly-added Datasets in ADBench. Since most of the public datasets are relatively small and simple, we introduce 7 more complex datasets from CV and NLP domains with more samples and richer features in ADBench (highlighted in Table 2 with blue). Pretrained models are applied to extract data embedding from NLP and CV datasets to access these more complex representations. For NLP datasets, we use BERT [47] pretrained on the BookCorpus and English Wikipedia to extract the embedding of the [CLS] token. For CV datasets, we use ResNet18 [38] pretrained on the ImageNet [22] to extract the embedding after the last average pooling layer. Following previous works [91, 92], we set one of the multi-classes as normal and downsample the remaining classes to 5% of the total instances as anomalies, and report the average results over all the respective classes.

Table 2: Data description of the 55 datasets included in ADBench; 7 newly added datasets from CV and NLP domain are highlighted in blue at the bottom of the table. To be consistent with the original datasets, we keep the original ratio of datasets, and some of them may be higher than 50%.

Data	# Samples	# Features	# Anomaly	% Anomaly	Category
abalone	4177	7	2081	49.82	Biology
ALOI	49534	27	1508	3.04	Image
annthyroid	7200	6	534	7.42	Healthcare
Arrhythmia	450	259	206	45.78	Healthcare
breastw	683	9	239	34.99	Healthcare
cardio	1831	21	176	9.61	Healthcare
Cardiotocography	2114	21	466	22.04	Healthcare
comm.and.crime	1994	101	993	49.80	Socio-economic
concrete	1030	8	515	50.00	Physical
cover	286048	10	2747	0.96	Botany
fault	1941	27	673	34.67	Physical
glass	214	7	9	4.21	Forensic
HeartDisease	270	13	120	44.44	Healthcare
Hepatitis	80	19	13	16.25	Healthcare
http	567498	3	2211	0.39	Web
imgseg	2310	18	990	42.86	Image
InternetAds	1966	1555	368	18.72	Image
Ionosphere	351	32	126	35.90	Oryctognosy
landsat	6435	36	1333	20.71	Astronautics
letter	1600	32	100	6.25	Image
Lymphography	148	18	6	4.05	Healthcare
magic.gamma	19020	10	6688	35.16	Physical
mammography	11183	6	260	2.32	Healthcare
mnist	7603	100	700	9.21	Image
musk	3062	166	97	3.17	Chemistry
optdigits	5216	64	150	2.88	Image
PageBlocks	5393	10	510	9.46	Document
Parkinson	195	22	147	75.38	Healthcare
pendigits	6870	16	156	2.27	Image
Pima	768	8	268	34.90	Healthcare
satellite	6435	36	2036	31.64	Astronautics
satimage-2	5803	36	71	1.22	Astronautics
shuttle	49097	9	3511	7.15	Astronautics
skin	245057	3	50859	20.75	Image
smtpt	95156	3	30	0.03	Web
SpamBase	4207	57	1679	39.91	Document
speech	3686	400	61	1.65	Linguistics
Stamps	340	9	31	9.12	Document
thyroid	3772	6	93	2.47	Healthcare
vertebral	240	6	30	12.50	Biology
vowels	1456	12	50	3.43	Linguistics
Waveform	3443	21	100	2.90	Physics
WBC	223	9	10	4.48	Healthcare
WDBC	367	30	10	2.72	Healthcare
Wilt	4819	5	257	5.33	Botany
wine	129	13	10	7.75	Chemistry
WPBC	198	33	47	23.74	Healthcare
yeast	1484	8	507	34.16	Biology
CIFAR10	5263	512	263	5.00	Image
FashionMNIST	6315	512	315	5.00	Image
SVHN	5208	512	260	5.00	Image
Agnews	10000	768	500	5.00	NLP
Amazon	10000	768	500	5.00	NLP
Imdb	10000	768	500	5.00	NLP
Yelp	10000	768	500	5.00	NLP

B.2 ADBench Algorithm List

We organize all the algorithms in ADBench into the following three categories and report their hyperparameter settings which mainly refer to the settings of their original papers or repositories (e.g., PyOD¹ and scikit-learn²).

(i) 14 unsupervised algorithms:

1. **Principal Component Analysis (PCA)** [97]. PCA is a linear dimensionality reduction using singular value decomposition of the data to project it to a lower dimensional space. When used for AD, it projects the data to the lower dimensional space and then uses the reconstruction errors as the anomaly scores. If not specified, the default hyperparameters in PyOD are used for the PCA (and the other unsupervised algorithms deployed by PyOD).
2. **One-class SVM (OCSVM)** [94]. OCSVM maximizes the margin between origin and the normal samples, and defines the decision boundary as the hyperplane that determines the margin.
3. **Local Outlier Factor (LOF)** [12]. LOF measures the local deviation of density of a given sample with respect to its neighbors.
4. **Clustering Based Local Outlier Factor (CBLOF)** [39]. CBLOF calculates the anomaly score by first assigning samples to clusters, and then using the distance among clusters as anomaly scores.
5. **Connectivity-Based Outlier Factor (COF)** [100]. COF uses the ratio of average chaining distance of data points and the average of average chaining distance of k -th nearest neighbor of the data point, as the anomaly score for observations.
6. **Histogram-based outlier detection (HBOS)** [29]. HBOS assumes the feature independence and calculates the degree of outlyingness by building histograms.
7. **K-Nearest Neighbors (KNN)** [85]. KNN views the anomaly score of the input instance as the distance to its k -th nearest neighbor.
8. **Subspace Outlier Detection (SOD)** [51]. SOD aims to detect outlier in varying subspaces of a high dimensional feature space.
9. **Copula Based Outlier Detector (COPOD)** [59]. COPOD is a hyperparameter-free, highly interpretable outlier detection algorithm based on empirical copula models.
10. **Empirical-Cumulative-distribution-based Outlier Detection (ECOD)** [60]. ECOD is a hyperparameter-free, highly interpretable outlier detection algorithm based on empirical CDF functions. Basically, it uses ECDF to estimate density of each feature independently, and assumes that outliers locate the tails of distribution.
11. **Deep Support Vector Data Description (DeepSVDD)** [91]. DeepSVDD trains a neural network while minimizing the volume of a hypersphere that encloses the network representations of the data, forcing the network to extract the common factors of variation.
12. **Deep Autoencoding Gaussian Mixture Model (DAGMM)** [124]. DAGMM utilizes a deep autoencoder to generate a low-dimensional representation and reconstruction error for each input data point, which is further fed into a Gaussian Mixture Model (GMM). We train the DAGMM for 200 epochs with 256 batch size, where the patience of early stopping is set to 50. The learning rate of Adam [49] optimizer is 0.0001 and is decayed once the number of epoch reaches 50. The latent dimension of DAGMM is set to 1 and the number of Gaussian mixture components is set to 4. The λ_1 and λ_2 for energy and covariance in the objective function is set to 0.1 and 0.005, respectively.
13. **Lightweight on-line detector of anomalies (LODA)** [80]. LODA is an ensemble method and is particularly useful in domains where a large number of samples need to be processed in real-time or in domains where the data stream is subject to concept drift and the detector needs to be updated on-line.
14. **Isolation Forest (IForest)** [63]. IForest isolates observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.

(ii) 7 semi-supervised algorithms:

1. **Semi-Supervised Anomaly Detection via Adversarial Training (GANomaly)** [5]. A GAN-based method that defines the reconstruction error of the input instance as the anomaly score. We replace the convolutional layer in the original GANomaly by the dense layer with tanh activation function 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 train the GANomaly for 50 epochs with 64 batch size, where the SGD [89] optimizer with 0.01 learning rate and 0.7 momentum is applied for both the generator and the discriminator.
2. **Deep Semi-supervised Anomaly Detection (DeepSAD)** [92]. A deep one-class method that improves the unsupervised DeepSVDD [91] by penalizing the inverse of the distances of anomaly representation such that anomalies must be mapped further away from the hypersphere center. The hyperparameter η in the loss function is set to 1.0, where DeepSAD is trained for 50 epochs with 128 batch size. Adam optimizer with 0.001 learning rate and 10^{-6} weight decay is applied for updating the network parameters. DeepSAD additionally employs an autoencoder for calculating the initial center of hypersphere, where the autoencoder

¹<https://pyod.readthedocs.io/en/latest/pyod.html>

²<https://scikit-learn.org/stable/>

is trained for 100 epochs with 128 batch size, and optimized by Adam optimizer with learning rate 0.001 and 10^{-6} weight decay.

3. **REPresentations for a random nEarest Neighbor distance based method (REPEN)** [74]. A neural network based model that leverages transformed low-dimensional representation for random distance-based detectors. The hidden size of REPEN is set to 20 and the margin of triplet loss is set to 1000. REPEN is trained for 30 epochs with 256 batch size, where the total number of steps (batches of samples) is set to 50. Adadelta [111] optimizer with 0.001 learning rate and 0.95 ρ is applied to update network parameters.
4. **Deviation Networks (DevNet)** [78]. A neural network based model that uses a prior probability to enforce statistically deviation score of input instances. The margin hyperparameter a in the deviation loss is set to 5. DevNet is trained for 50 epochs with 512 batch size, where the total number of steps is set to 20. RMSprop [89] optimizer with 0.001 learning rate and 0.95 ρ is applied to update network parameters.
5. **Pairwise Relation prediction-based ordinal regression Network (PReNet)** [77]. A neural network based model that defines a two-stream ordinal regression to learn the relation of instance pairs. The score targets of {unlabeled, unlabeled}, {labeled, unlabeled} and {labeled, labeled} sample pairs are set to 0, 4 and 8, respectively. PReNet is trained for 50 epochs with 512 batch size, where the total number of steps is set to 20. RMSprop optimizer with learning rate 0.001 and 0.01 weight decay is applied to update network parameters.
6. **Feature Encoding With Autoencoders for Weakly Supervised Anomaly Detection (FEAWAD)** [122]. A neural network based model that incorporates the network architecture of DAGMM [124] with the deviation loss of DevNet [78]. FEAWAD is trained for 30 epochs with 512 batch size, where the total number of steps is set to 20. Adam optimizer with 0.0001 learning rate is applied to update network parameters.
7. **Extreme Gradient Boosting Outlier Detection (XGBOD)** [115]. XGBOD first uses the passed in unsupervised outlier detectors to extract richer representation of the data and then concatenates the newly generated features to the original feature for constructing the augmented feature space. An XGBoost classifier is then applied on this augmented feature space. We use the default hyperparameters in PyOD.

(iii) *9 supervised algorithms:*

1. **Naive Bayes (NB)** [9]. NB methods are based on applying Bayes’ theorem with the “naive” assumption of conditional independence between every pair of features given the value of the class variable. We use the Gaussian NB in ADBench.
2. **Support Vector Machine (SVM)** [19]. SVM is effective in high dimensional spaces and could be still effective in cases where the number of dimensions is greater than the number of samples. We use the default hyperparameters in scikit-learn for SVM (and for the following MLP and RF).
3. **Multi-layer Perceptron (MLP)** [88]. MLP uses the binary cross entropy loss to update network parameters.
4. **Random Forest (RF)** [11]. RF is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting.
5. **eXtreme Gradient Boosting (XGBoost)** [17]. XGBoost is an optimized distributed gradient boosting method designed to be highly efficient, flexible and portable. We use the default hyperparameter settings in the XGBoost official repository¹.
6. **Highly Efficient Gradient Boosting Decision Tree (LightGBM)** [46]. LightGBM is a gradient boosting framework that uses tree based learning algorithms with faster training speed, higher efficiency, lower memory usage and better accuracy. The default hyperparameter settings in the LightGBM official repository² are used.
7. **Categorical Boosting (CatBoost)** [81]. CatBoost is a fast, scalable, high performance gradient boosting on decision trees. CatBoost uses the default hyperparameter settings in its official repository³.
8. **Residual Nets (ResNet)** [33]. This method introduces a ResNet-like architecture [38] for tabular based data. ResNet is trained for 100 epochs with 64 batch size. AdamW [66] optimizer with 0.001 learning rate is applied to update network parameters.
9. **Feature Tokenizer + Transformer (FTTransformer)** [33]. FTTransformer is an effective adaptation of the Transformer architecture [103] for tabular data. FTTransformer is trained for 100 epochs with 64 batch size. AdamW optimizer with 0.0001 learning rate and 10^{-5} weight decay is applied to update network parameters.

B.3 Open-source Release

As mentioned before, the full experiment code, datasets, and the examples of benchmarking new algorithms are available at <https://github.com/Minqi824/ADBench>. We specify the key environment setting of using ADBench, e.g., `scikit-learn==0.20.3`, `pyod==0.9.8`, etc. With our interactive example in Jupyter notebooks, one may compare newly proposed AD algorithm easily.

¹<https://xgboost.readthedocs.io/en/stable/parameter.html>

²<https://lightgbm.readthedocs.io/en/latest/Parameters.html>

³<https://catboost.ai/en/docs/references/training-parameters/>

C Details on Experiment Setting

We provide additional details on experiment setting to §4.1 in this section.

General Experimental Settings. Although unsupervised AD algorithms are primarily designed for the transductive setting (i.e., outputting the anomaly scores on the input data only other than making predictions on newcoming data), we adapt all the algorithms for the inductive setting to make prediction on the newcoming data, which is helpful in applications and also common in popular AD library PyOD [118], TODS [52, 53], and PyGOD [64]. Thus, we use 70% data for training and the remaining 30% as testing set. We use stratified sampling to keep the anomaly ratio consistent. We repeat each experiment 3 times and report the average. The 7 complex CV and NLP datasets are mainly considered for evaluating algorithm performance on the real-world datasets and are not included in the experiments of different types of anomalies and algorithm robustness, since such high-dimensional data could make it hard to generate synthetic anomalies (e.g., the Vine Copula is computationally expensive for fitting such high-dimensional data), or introduce too much noises in input data (e.g., the noise ratio of irrelevant features 50% would lead to 384 noise features in the 768 input dimensions of NLP data).

Hyperparameter Settings. For all the algorithms in ADBench, we use their default hyperparameter (HP) settings in the original paper for fair comparison. Specific values can be found in Appx.B.2 and our codebase¹. It is also acknowledged that it is possible to use a small hold-out data for hyperparameter tuning for semi- and fully-supervised methods [98], while we do not consider this setting in this work.

Extensive Experiments. In total ADBench conducts 93,654 experiments, where each denotes one algorithm’s result on a dataset under a specific setting. More specifically, we have 20,790 experiments in §4.2:

- Unsupervised methods on real-world datasets {14 algorithms, 55 datasets, 3 repeat times} leads to 2,310 experiments.
- Semi- and fully-supervised on real-world datasets {16 algorithms, 55 datasets, 3 repeat times, 7 settings of labeled anomalies} leads to 18,480 experiments.

Additionally, we have 18,144 experiments for understanding the algorithm performances under four types of anomalies in §4.3:

- Unsupervised methods on real-world datasets {14 algorithms, 48 datasets, 3 repeat times} leads to 2,016 experiments.
- Semi- and fully-supervised on real-world datasets {16 algorithms, 48 datasets, 3 repeat times, 7 settings of labeled anomalies} leads to 16,128 experiments.

Finally, we have 54,720 experiments for evaluating the algorithm robustness under three settings of data noises and corruptions in §4.4:

- For duplicated anomalies and irrelevant features {30 algorithms, 48 datasets, 3 repeat times, 5 settings of data noises, 2 scenarios} leads to 43,200 experiments.
- For annotation errors {16 algorithms, 48 datasets, 3 repeat times, 5 settings of data noises} leads to 11,520 experiments.

Computational Resources. Classical anomaly detection models are run on an Intel i7-8700 @3.20 GHz, 16GB RAM, 12-core workstation. For deep learning models (especially for ResNet and FTTransformer), we run the experiments on an NVIDIA Tesla V100 GPU accelerators. The model runtime on real-world datasets is reported in Appx. §D.1.

D Additional Experiment Results

D.1 Additional Results for Overall Model Performance on Real-world Datasets in §4.2

In addition to the AUCROC results presented in §4.2, we also show the AUCPR results of model performance on 55 real-world datasets in Fig. 8, where the corresponding conclusions are similar to that of AUCROC results. There is still no statistically superior solution for unsupervised methods regarding AUCPR. Semi-supervised methods perform better than supervised methods when only limited label data is available, say the labeled anomalies γ_l is less than 5%. Besides, we show that the semi-supervised GANomaly, which learns an intermediate representation of the normal data, performs worse than those anomaly-informed model leveraging labeled anomalies (see Fig. 8(b)). This conclusion verifies that merely capturing the normal behaviors is not enough for detecting the underlying anomalies, where the lack of knowledge about the true anomalies would lead to high false positives / negatives [75, 77, 78].

¹ADBench repo: <https://github.com/Minqi824/ADBench>

Fig. 9 and 10 show the boxplots of AUCROC and AUCPR of 30 algorithms on the 55 real-world datasets. These results validate the no-free-lunch theorem, where there is no model that is both the best and the most stable performer. For example, DeepSVDD and RF are the most stable detectors among un- and fully-supervised methods, respectively, but they are inferior to most of the other algorithms. Besides, IForest and CatB(oost) can be regarded as two very competitive methods among un- and fully-supervised methods, respectively, but their variances of model performance are relatively large compared to the other methods.

Fig. 11 reports the runtime of each compared algorithm. Generally deep learning methods like ResNet and FTTransformer and ensemble methods like XGBOD need more time for model training and convergence, where some shallow unsupervised methods like HBOS, COPOD and ECOD usually run faster.

Additionally, we also present the full results in tables in §D.4.

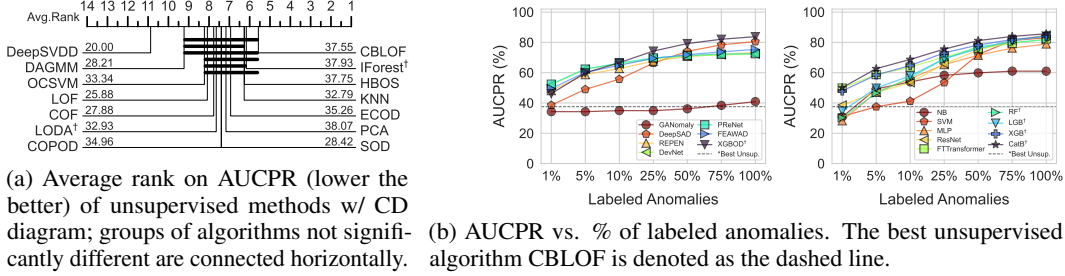


Figure 8: AD model’s AUCPR on 55 real-world datasets. Generally the AUCPR results are consistent with the AUCROC results in §4.2. (a) shows that no unsupervised algorithm can statistically outperform. (b) shows the AUCPR of semi- and supervised methods under varying ratio of labeled anomalies γ_l . The semi-supervised methods leverage the labels more efficiently w/ small γ_l .

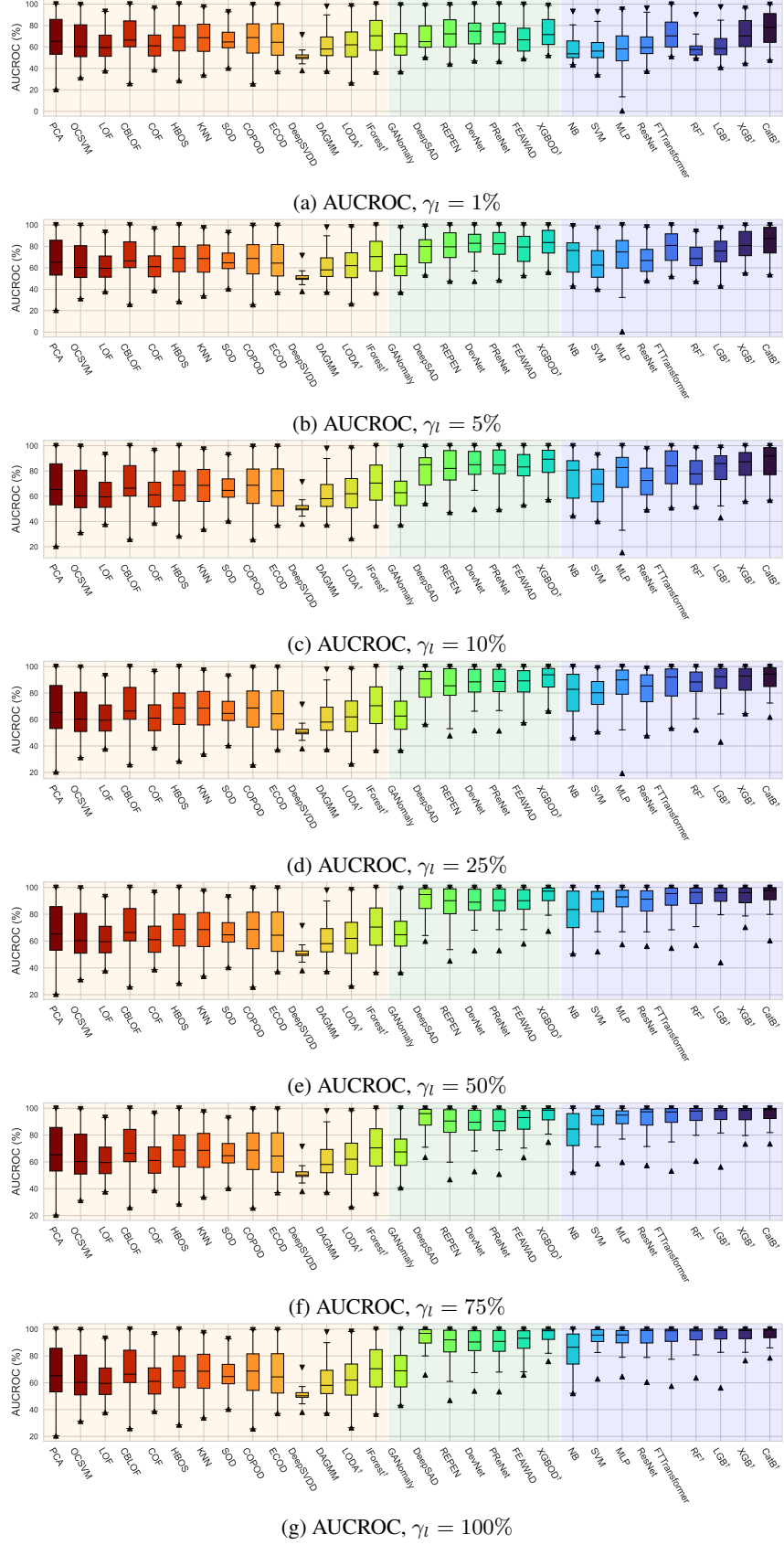


Figure 9: Boxplot of AUCROC. We denote unsupervised methods in , semi-supervised methods in , and supervised methods in . Consistent with the CD diagrams, we notice that none of the unsupervised methods visually outperform.

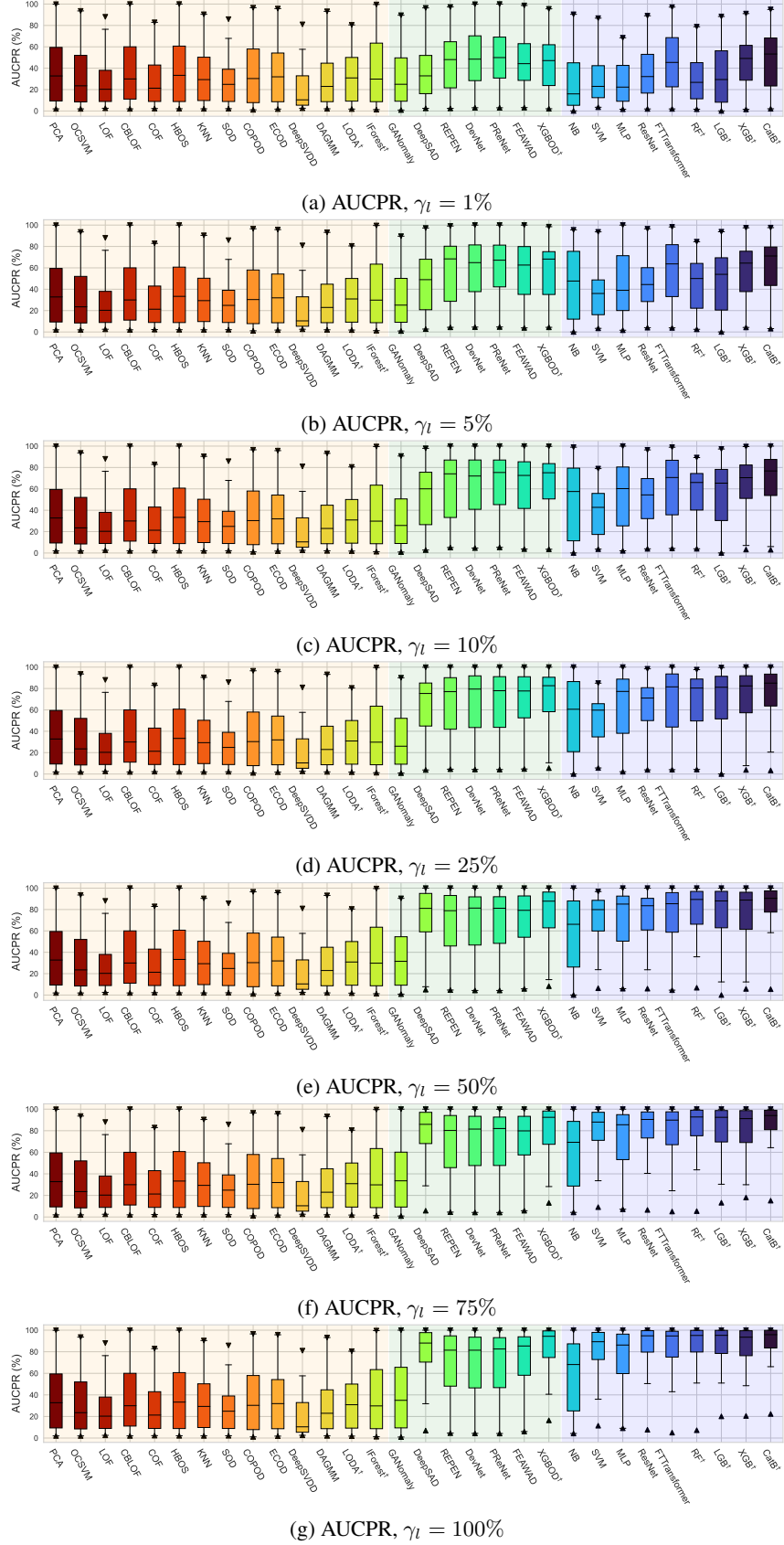


Figure 10: Boxplot of AUCPR. We denote unsupervised methods in ■, semi-supervised methods in ■, and supervised methods in ■. Consistent with the CD diagrams, we notice that none of the unsupervised methods visually outperform.

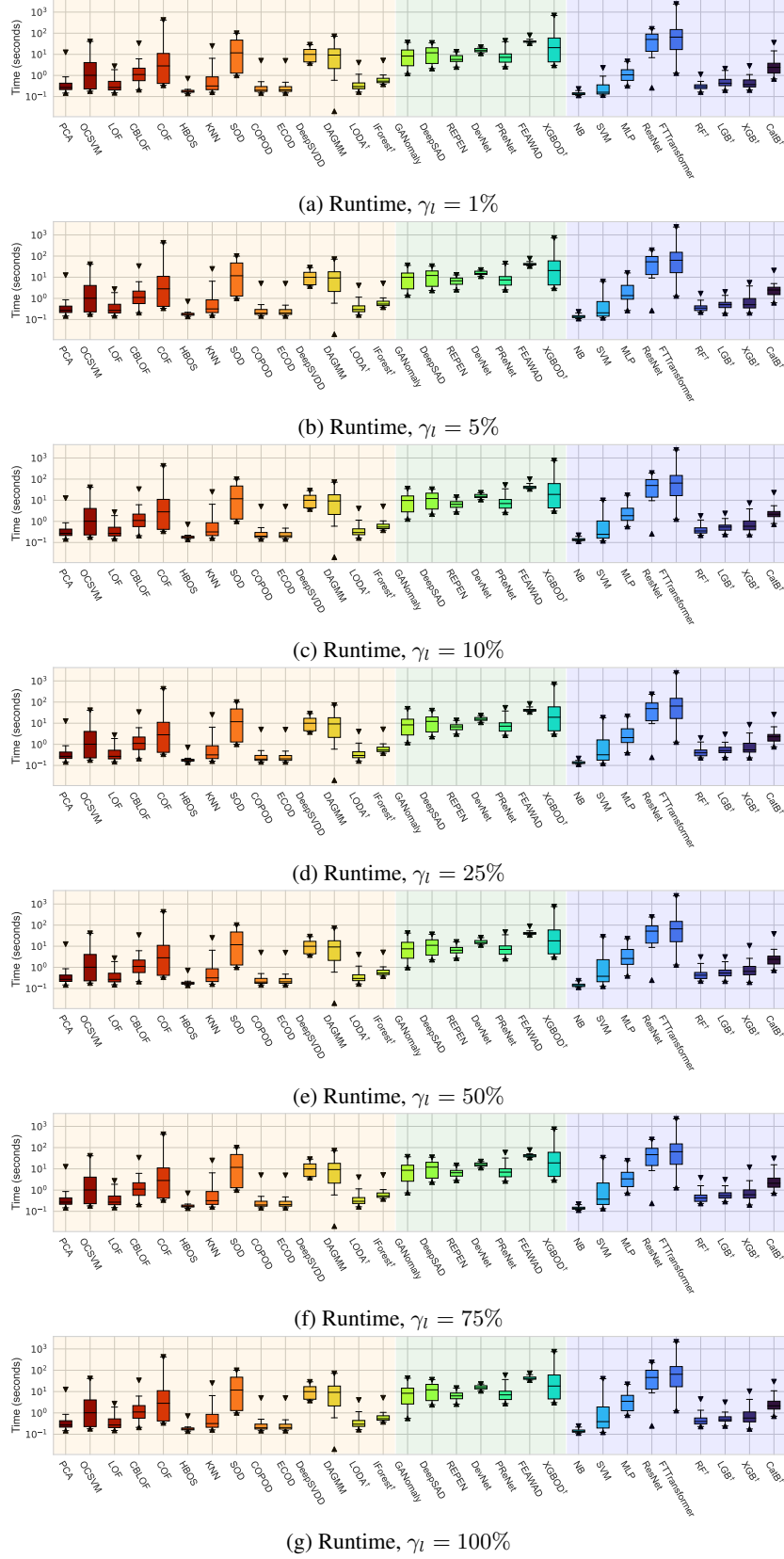


Figure 11: Runtime of included algorithms. We denote unsupervised methods in , semi-supervised methods in , and supervised methods in . We find HBOS, COPOD, ECOD, and NB are fastest as they treat each feature independently. In contrast, more complex feature representation methods like XGBOD, ResNet, and FTTTransformer appear to be computationally heavy.

D.2 Additional Results for Different Types of Anomalies §4.3

We additionally show the AUCPR results for model performance on different types of anomalies in Fig. 12 and Fig. 13, which are consistent with the conclusions drawn in §4.3, i.e., the unsupervised methods are significantly better if their model assumptions conform to the underlying anomaly types. Moreover, the prior knowledge of anomaly types can be more important than that of label information, where those label-informed algorithms generally underperform the best unsupervised methods for local, global and dependency anomalies.

We want to note that XGBOD can be regarded as an exception to the above observations, which is comparable or even outperforms the best unsupervised model when more labeled anomalies are available. Recall that XGBOD employs the stacking ensemble method [105], where heterogeneous unsupervised methods are integrated with the supervised model XGBoost, therefore XGBOD is more adaptable to different data assumptions while effectively leveraging the label information. This validates the conclusion that such ensemble learning techniques should be considered in the future research direction.

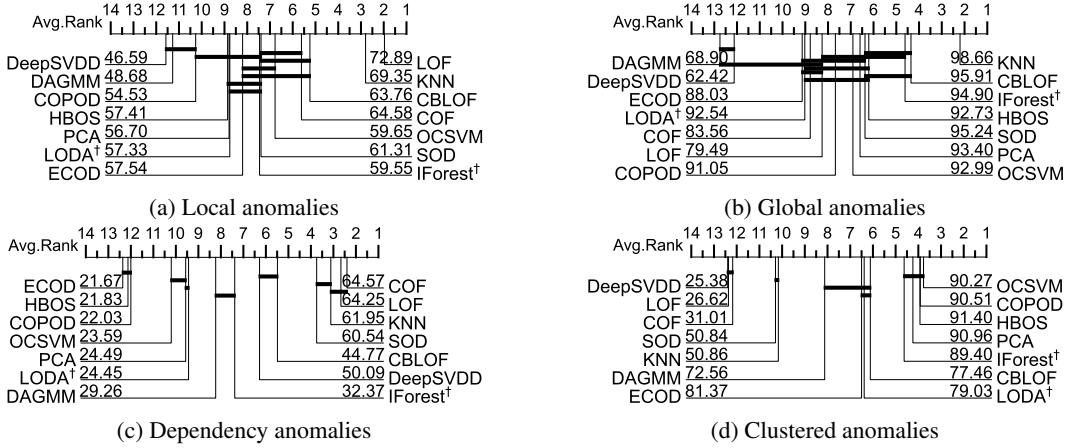


Figure 12: AUCPR CD Diagram of unsupervised methods on different types of anomalies. The unsupervised methods perform well when their assumptions conform to the anomaly types.

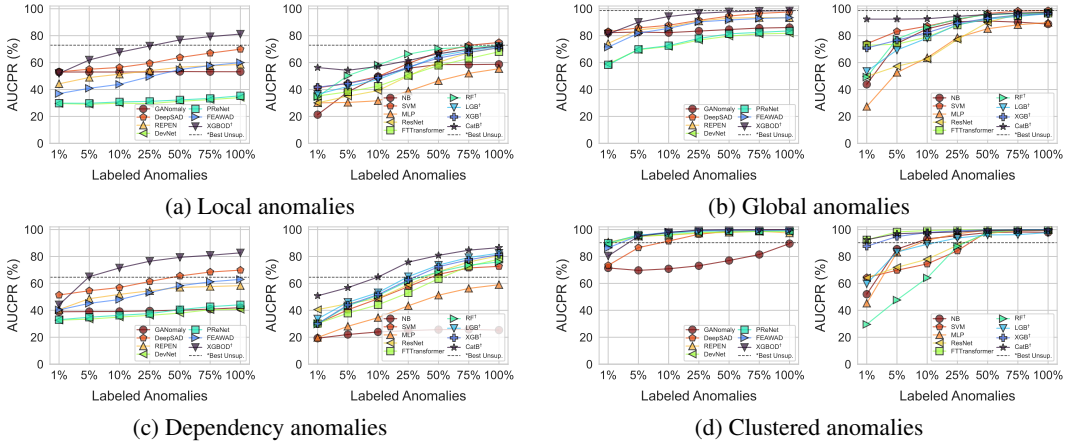


Figure 13: Semi- (left of each subfigure) and supervised (right) algorithms' performance on different types of anomalies with varying level of labeled anomalies for AUCPR performance. Surprisingly, these label-informed algorithms are *inferior* to the best unsupervised method except for the clustered anomalies.

D.3 Additional Results for Algorithm Robustness in §4.4

In addition to the primary results shown in §4.4, we provide the AUCPR results for algorithm robustness in Fig. 14. The AUCPR results confirm the robustness of supervised methods for irrelevant features. Besides, both semi- and fully-supervised methods are robust to minor annotation errors, say the annotation errors are less than 10%.

One thing to note is we observe AUCPR performance improves under the setting of duplicated anomalies (see Fig. 14 (a)-(c)). This is expected as AUCPR puts more emphasis on the positive classes (i.e., anomalies), and more duplicated anomalies favors this metric. Since this observation is consistently true for both unsupervised methods and label-informed methods, it would not largely impact our selection of algorithms. However, if we care both anomaly and normal classes equally, the results on AUCROC in §4.4 still stands—unsupervised methods are more susceptible to duplicate anomalies.

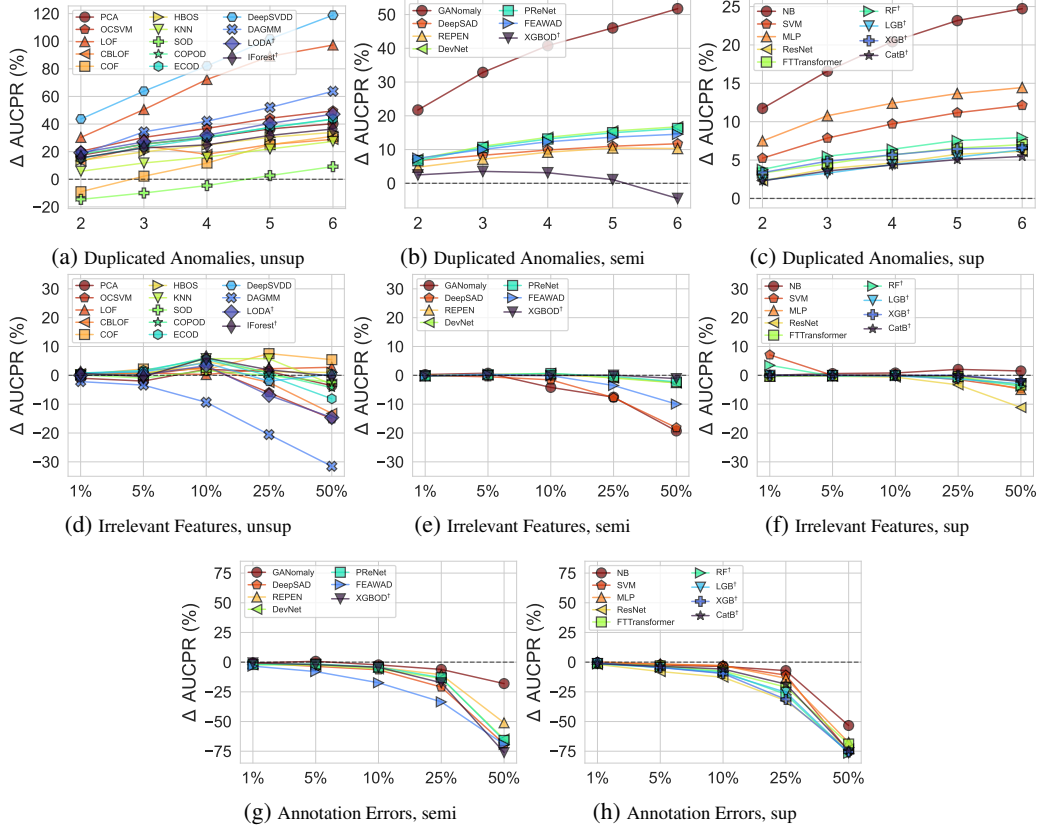


Figure 14: Algorithm performance change under noisy and corrupted data (i.e., duplicated anomalies for (a)-(c), irrelevant features for (d)-(f), and annotation errors for (g) and (h)). y-axis denotes the % of performance change (ΔAUCPR) and its range remains consistent across different algorithms. The results reveal the usage of label information in defending irrelevant features, and the robustness of label-informed methods to the minor annotation errors. Un-, semi-, and fully-supervised methods are denoted as *unsup*, *semi*, and *sup*, respectively. The results are mostly consistent with the observations in Fig. 7 (§4.4).

D.4 Full Performance Tables on Real-world Datasets (in addition to §4.2 and Appendix D.1)

In the following tables, we first present the AUCROC and AUCPR for all unsupervised methods, and then show the label-informed methods’ performance at different levels of labeled anomaly ratio (i.e., $\gamma_l = \{1\%, \dots, 100\%\}$). We would expect these results are useful in constructing unsupervised anomaly detection model selection methods like MetaOD [119], where the historical algorithm performance table serves as great sources for building strong meta-learning methods.

Table 3: AUCROC of 14 unsupervised algorithms on 55 real-world datasets. We show the performance rank in parenthesis (lower the better).

Datasets	PCA	OCSVM	LOF	CBLOF	COF	HBOS	KNN	SOD	COPOD	ECOD	Deep SVDD	DA GMM	LODA	IForest
abalone	49.08(9)	48.66(10)	57.94(4)	59.18(3)	57.82(6)	47.91(13)	77.45(1)	69.81(2)	55.6(7)	42.82(14)	49.82(8)	48.62(1)	57.86(5)	48.64(11)
ALOI	56.65(6)	55.85(8)	66.63(1)	55.22(9)	64.68(2)	52.63(12)	61.47(3)	61.09(4)	53.75(10)	56.6(7)	52.99(11)	51.96(13)	51.33(14)	56.66(5)
anthyroid	66.24(7)	57.23(11)	70.2(6)	62.26(9)	65.92(8)	60.15(10)	71.69(5)	77.38(3)	76.8(4)	78.03(2)	49.82(13)	56.53(12)	41.02(14)	82.01(1)
Arrhythmia	73.89(7)	73.55(3)	72.99(9)	73.02(8)	63.31(12)	74.51(4)	74.46(5)	64.63(11)	76.5(2)	77.88(1)	56.4(14)	62.38(13)	70.91(10)	74.38(6)
breastw	95.13(8)	80.3(10)	40.61(12)	96.81(7)	38.84(13)	98.94(3)	97.01(6)	93.97(9)	99.68(1)	99.17(2)	66.97(11)	N/A(N/A)	98.49(4)	98.32(5)
cardio	95.55(1)	93.91(3)	66.33(13)	89.93(7)	71.41(12)	84.67(8)	76.64(9)	73.25(11)	92.35(5)	94.44(2)	48.99(14)	75.01(10)	90.34(6)	93.19(4)
Cardiotocography	74.67(2)	77.86(1)	59.51(10)	64.54(7)	53.77(12)	60.86(9)	56.23(11)	51.69(13)	67.02(6)	68.92(4)	51.2(14)	62.01(8)	73.65(3)	67.57(5)
command.crim	58.2(9)	54.28(12)	52.57(13)	62.68(6)	54.97(11)	67.34(1)	66.41(2)	61.02(7)	65.89(3)	62.69(5)	52.19(14)	54.90(10)	59.32(8)	63.36(4)
concrete	55.58(13)	51.68(14)	64.22(5)	62.37(6)	61.72(7)	70.74(1)	61.38(8)	57.87(11)	67.1(3)	60.18(9)	57.19(12)	59.06(10)	68.7(2)	65.64(4)
cover	93.73(1)	92.62(3)	84.58(10)	89.3(6)	76.91(12)	80.24(11)	85.97(9)	74.46(13)	88.64(7)	93.42(2)	48.91(14)	89.89(5)	92.34(4)	86.74(8)
fault	46.02(10)	47.69(9)	58.93(5)	64.06(3)	62.1(4)	51.28(7)	72.98(1)	68.11(2)	43.88(12)	43.41(13)	48.7(8)	45.86(11)	41.71(14)	57.02(6)
glass	66.29(12)	35.36(14)	69.21(1)	82.94(1)	72.24(10)	77.23(3)	82.29(2)	73.36(7)	72.43(9)	70.71(6)	66.64(8)	76.99(5)	73.13(8)	77.13(4)
HeartDisease	60.33(4)	54.41(9)	51.25(11)	64.59(3)	51.45(10)	55.19(1)	55.86(8)	56.93(7)	70.2(5)	73.82(1)	50.13(14)	73.41(12)	78.42(9)	84.5(6)
Hepatitis	75.95(4)	67.75(7)	38.02(14)	66.4(8)	41.45(13)	79.85(2)	52.76(11)	68.17(6)	82.05(1)	79.67(3)	45.83(12)	54.8(10)	64.87(9)	69.75(5)
http	99.72(2)	99.59(4)	27.46(11)	99.6(3)	88.78(8)	99.53(5)	3.37(13)	78.04(9)	99.29(6)	98.1(7)	49.96(10)	N/A(N/A)	12.48(12)	99.96(1)
imgseg	58.4(7)	59.68(5)	47.47(12)	50.62(9)	50.05(11)	64.14(3)	43.94(14)	45.29(13)	71.39(1)	68.69(2)	50.19(10)	51.88(8)	61.88(4)	58.81(6)
InternetAds	61.67(12)	68.28(4)	65.83(9)	70.58(1)	63.79(10)	68.03(5)	69.99(2)	61.85(11)	67.05(7)	67.1(6)	66.64(8)	N/A(N/A)	55.38(13)	69.01(3)
Ionosphere	79.19(8)	75.92(10)	90.59(2)	90.72(1)	86.76(4)	62.49(13)	88.26(3)	86.41(5)	79.34(7)	75.59(11)	50.13(14)	73.41(12)	78.42(9)	84.5(6)
landsat	35.76(14)	36.15(13)	53.9(6)	63.55(1)	53.5(7)	55.14(5)	57.95(3)	59.54(2)	41.55(11)	56.61(4)	49.93(8)	43.92(10)	38.17(12)	47.64(9)
letter	50.29(11)	46.18(13)	84.49(2)	75.62(5)	80.03(4)	59.74(7)	86.19(1)	84.09(3)	54.32(8)	50.76(9)	38.02(14)	50.42(10)	50.24(12)	61.07(6)
Lymphography	89.82(2)	99.54(4)	89.86(9)	99.83(1)	90.85(8)	99.49(6)	55.91(13)	72.49(11)	99.48(7)	99.52(5)	49.83(14)	72.11(12)	85.55(10)	99.81(3)
magic.gamma	67.22(9)	60.65(12)	68.51(6)	75.13(3)	66.64(10)	70.86(5)	82.38(1)	75.4(2)	68.33(7)	64.36(11)	49.85(14)	58.32(8)	68.02(8)	73.25(4)
mammography	88.72(3)	84.95(6)	74.39(12)	83.74(9)	77.53(11)	86.27(5)	84.53(7)	81.51(10)	90.69(2)	90.75(1)	52.9(13)	N/A(N/A)	83.91(8)	86.39(4)
mnist	85.29(1)	82.95(3)	67.13(1)	79.45(6)	60.42(12)	80.58(5)	60.11(13)	77.74(7)	84.6(2)	51.51(14)	84.6(2)	51.51(14)	67.23(10)	72.27(8)
musk	100(1)	80.58(8)	41.18(13)	100(1)	38.69(14)	100(1)	69.89(11)	74.09(10)	94.2(7)	95.11(5)	50.59(12)	76.85(9)	95.11(5)	99.99(4)
optdigits	51.65(12)	54(11)	56.1(9)	87.51(1)	49.15(13)	81.63(2)	41.73(14)	58.92(8)	68.71(4)	61.04(7)	54.24(10)	62.57(5)	61.74(6)	70.92(3)
PageBlocks	90.64(2)	88.76(5)	75.9(12)	85.04(7)	72.65(13)	80.58(10)	81.94(9)	77.75(11)	88.05(6)	90.92(1)	51.7(14)	89.61(3)	83.34(8)	89.57(4)
Parkinson	35.73(13)	32.91(14)	55.89(6)	60.56(3)	59.14(4)	68.79(1)	52.33(8)	65.45(2)	54.25(7)	37.11(12)	58.41(5)	40.55(11)	49.02(10)	50.5(9)
pendigits	93.37(2)	93.75(2)	47.99(13)	90.4(7)	45.07(14)	93.04(4)	72.95(9)	66.29(10)	90.68(6)	91.22(5)	48.42(12)	47.62(11)	49.1(8)	94.76(1)
Pima	70.77(5)	66.92(7)	65.71(9)	71.42(3)	61.05(11)	71.07(4)	73.43(1)	61.25(10)	69.1(6)	51.54(13)	46.19(14)	55.92(12)	65.93(8)	72.87(2)
satellite	59.62(10)	59.02(11)	55.88(12)	71.32(3)	54.74(13)	74.8(2)	65.18(5)	63.96(6)	63.2(7)	75.06(1)	50.6(14)	62.33(8)	61.98(9)	70.34(4)
satimage-2	97.62(4)	97.35(6)	47.36(14)	99.84(1)	56.7(12)	97.65(3)	92.6(10)	83.08(11)	97.21(7)	97.11(8)	50.83(13)	96.29(9)	97.56(5)	99.16(2)
shuttle	98.62(5)	97.4(7)	57.11(12)	83.48(8)	51.72(13)	98.63(4)	69.64(9)	69.51(10)	99.35(3)	99.4(2)	30.81(14)	97.92(6)	60.95(11)	99.56(1)
skin	45.26(1)	49.45(7)	60.15(9)	61.45(8)	47.14(1)	85.56(13)	73.79(1)	40.32(6)	25.64(14)	37.51(12)	58.02(6)	N/A(N/A)	45.75(10)	63.21(3)
smtp	88.41(3)	80.7(4)	71.84(9)	79.68(5)	79.6(6)	70.52(11)	89.62(2)	59.85(13)	79.09(7)	71.86(8)	58.33(14)	71.32(10)	67.43(12)	89.73(1)
SpamBase	54.66(6)	52.47(8)	43.33(11)	54.97(5)	40.96(13)	64.74(4)	53.35(7)	52.35(9)	70.09(1)	66.89(2)	50.42(10)	N/A(N/A)	41.99(12)	64.76(3)
speech	50.79(9)	50.19(13)	52.48(6)	50.58(12)	55.97(2)	50.59(11)	51.03(8)	55.86(3)	52.89(4)	51.58(7)	56.58(1)	52.75(5)	49.84(14)	50.74(10)
Stamps	91.47(2)	83.86(8)	51.26(13)	68.18(11)	53.81(12)	90.73(5)	68.61(10)	73.26(9)	93.4(1)	91.41(3)	45.1(14)	88.86(6)	87.18(7)	91.21(4)
thyroid	96.34(3)	87.92(10)	86.86(1)	94.73(6)	90.87(9)	95.62(5)	95.93(4)	92.81(8)	94.3(7)	97.78(2)	51.2(14)	79.75(12)	74.31(3)	98.3(1)
vertebral	31.41(5)	37.99(7)	49.29(3)	41.41(5)	48.71(4)	25.56(13)	35.61(14)	37.51(12)	49.43(2)	37.51(12)	49.43(2)	37.51(12)	49.43(2)	37.51(12)
vowels	65.29(9)	61.59(10)	93.12(3)	89.92(5)	94.04(2)	72.21(7)	97.26(1)	92.65(4)	53.15(12)	45.81(13)	44.25(14)	60.58(1)	70.36(8)	73.94(6)
Waveform	65.48(10)	56.29(12)	73.32(3)	72.42(6)	72.56(5)	68.77(8)	73.78(2)	68.57(9)	75.03(1)	73.25(4)	49.67(13)	49.35(14)	60.13(11)	71.47(7)
WBC	99.02(7)	99.03(4)	54.17(13)	99.46(1)	60.9(11)	98.72(6)	90.56(10)	94.6(9)	99.11(2)	99.11(2)	59.15(12)	N/A(N/A)	96.91(8)	99.01(5)
WDBC	99.05(4)	98.86(6)	89.12(9)	99.32(3)	96.26(9)	99.5(1)	91.72(11)	91.91(10)	99.42(2)	97.2(8)	71.42(14)	76.67(13)	98.26(7)	98.95(5)
Wilt	55.12(6)	62.7(3)	40.79(12)	44.25(9)	41.87(11)	58.94(4)	37.59(14)	39.16(13)	66.62(1)	64.96(2)	42.7(10)	45.7(8)	58.07(5)	52.71(7)
wine	84.37(4)	73.07(6)	37.74(13)	25.86(14)	44.44(12)	91.36(1)	44.98(11)	46.11(10)	88.65(3)	71.34(7)	50.35(9)	61.7(8)	90.12(2)	80.37(5)
WPBC	46.01(10)	45.53(12)	41.41(14)	44.77(13)	45.88(11)	51.24(2)	46.59(9)	51.14(3)	49.34(4)	46.83(7)	55.25(1)	47.8(6)	49.31(5)	46.63(8)
yeast	41.15(7)	41(9)	45.31(2)	44.85(3)	44.48(5)	39.64(10)	39.66(12)	42.46(6)	36.99(14)	39.61(11)	48.1(1)	41.11(8)	44.58(4)	37.76(13)
FashionMNIST	86.09(3)	85.24(4)	67.57(12)	88.17(1)	71.44(11)	78.68(10)	86.6(2)	81.73(7)	81.07(8)	83.63(6)	64.02(14)	67.29(13)	80.28(9)	84.89(5)
CFAR10	63.87(6)	63.76(7)	68.57(1)	64.23(4)	64.7(3)	57.51(3)	64.75(2)	64.22(5)	58.64(11)	61.04(10)	56.94(14)	58.08(12)	62.34(8)	61.28(9)
SVHN	95.16(3)	94.07(5)	64.51(1)	60.34(7)	63.47(2)	96.08(1)	96.08(1)	92.63(3)	91.09(4)	97.57(1)	58.27(9)	53.81(14)	88.26(10)	94.46(1)
Agnews	54.78(5)	54.34(9)	71.8(1)	60.02(5)	68.97(2)	53.87(10)	64.11(3)	62.81(4)	52.98(12)	53.04(11)	45.61(14)	52.02(13)	55.47(7)	56.74(6)
Amazon	55.06(10)	54.14(12)	56.11(9)	57.36(3)	56.96(4)	56.52(7)	60.03(2)	60.05(1)	56.94(5)	56.79(6)	51.04(14)	53.58(13)	54.2(11)	56.13(8)
lmdb	47.06(12)	46.07(14)	48.71(8)	49.35(5)	49.64(4)	49.1(6)	47.83(10)	49.86(3)	50.68(2)	50.73(1)	48.19(9)	47.67(11)	46.43(13)	49.09(7)
Yelp	60.71(11)	60.28(12)	67.09(3)	64.9(5)	66.11(4)	61.85(9)	69.84(1)	67.74(2)	62.56(7)	62.15(8)	49.67(14)	56.28(13)	61.36(10)	62.53(6)

Table 4: AUCPR of 14 unsupervised algorithms on 55 real-world datasets. We show the performance rank in parenthesis (lower the better).

Datasets	PCA	OCSVM	LOF	CBLOF	COF	HBOS	KNN	SOD	COPOD	ECOD	Deep SVDD	DA GMM	LODA	IForest
abalone	50.41(10)	50.75(9)	54.6(7)	56.3(5)	55.47(6)	53.58(8)	74.03(1)	67.63(2)	56.79(4)	42.9(14)	49.74(12)	49.04(13)	59.43(3)	50(11)
ALOI	4.17(10)	5.02(5)	8.08(1)	4.46(8)	6.85(2)	3.99(13)	6.02(3)	5.97(4)	3.62(14)	3.9(11)	4.59(6)	4.33(9)	4.53(7)	3.9(11)
anthyroid	16.12(7)	10.37(11)	15.71(8)	16.69(10)	14.39(9)	16.99(4)	16.74(5)	18.84(3)	16.58(6)	24.65(2)	7.84(13)	9.64(12)	7.06(14)	30.47(1)
Arrhythmia	70.21(10)	74.67(4)	71.57(8)	72.61(7)	63.51(11)	76.43(5)	74.5(6)	61.68(13)	76.87(2)	75.52(14)	57.52(14)	63.19(12)	70.24(9)	74.63(5)
breastw	95.16(1)	82.71(6)	28.55(12)	91.54(8)	27.6(13)	97.71(3)	92.19(7)	84.88(9)	99.4(1)	98.54(2)	61.45(11)	89.14(4)	98.04(5)	99.99(4)
cardio	66.06(2)	62.89(3)	23.79(13)	61.95(4)	28.67(11)	52.1(8)	40.72(9)	28.54(12)	60.42(5)	68.59(1)	21.19(14)	28.92(10)	53.41(7)	59.95(6)
Cardiotocography	47.95(3)	52.61(1)	30.66(10)	45.44(4)	28.21(12)	38.28(8)	34.79(9)	27.99(13)	40.46(7)	43.57(5)	24.97(14)	30.61(11)	48(2)	41.47(6)
command.crim	61.28(7)	58.47(10)	53.02(14)	61.71(6)	55.19(11)	68.42(1)	66.32(2)	61.25(8)	65.03(3)	62.18(5)	53.04(13)	53.72(12)	61.01(9)	64.07(4)
concrete	60.72(9)	57.19(12)	60.65(10)	64.27(3)	62.74(5)	67.46(2)	61.47(6)	56.94(13)	61.11(7)	56.19(14)	57.68(11)	61.05(8)	68.4(1)	63.83(4)
cover	9.86(1)	11.41(4)	8.12(8)	5.83(11)	4.1(2)	6.83(9)	6.16(10)	3.88(13)	11.37(5)	15.63(2)	2.55(14)	27.59(1)	13.06(3)	8.85(7)
fault	32.76(11)	38.44(6)	38.38(7)	43.98(3)	15.56(4)	36.47(8)	54.45(1)	48.01(2)	30.54(14)	30.82(13)	35.07(9)	33.48(10)	31.06(12)	41.09(5)
glass	10.05(14)	8.02(13)	2.11(13)	13.84(6)	11.81(9)	11.82(8)	20.26(2)	18.73(4)	9.78(12)	3.75(14)	24.58(1)	33.7(7)	10.99(10)	10.99(10)
Hepatitis	51.44(4)	45.48(7)	44.08(12)	48.08(3)	71.3(1)	59.04(2)	53.56(7)	46.36(9)	50.56(10)	40.42(13)	35.64(12)	32.66(11)	32.66(11)	32.66(11)
Hepatitis	36.65(4)	29.44(7)	16.69(13)	31.54(5)	14.39(13)	37.73(3)	21.95(11)	28.94(1)	41.5(1)	37.83(2)	19.01(12)	22.96(10)	33.96(7)	26.25(8)
http	56.43(2)	46.86(4)	3.82(10)	47.53(3)	9.57(8)	44.79(5)	0.71(11)	8.32(9)	35.19(6)	16.61(7)	0.38(13)	N/A(N/A)	0.67(12)	90.83(1)
imgsz	55.12(6)	62(73)	40.79(12)	44.25(9)	41.87(11)	58.94(4)	37.59(14)	39.16(13)	66.63(1)	64.96(2)	42.7(10)	58.10(8)	58.10(8)	52.71(7)
InternetAds	32.55(11)	54.68(2)	40.49(8)	58.13(1)	38.67(9)	53.97(3)	43.23(7)	27.69(12)	50.97(5)	51.07(4)	33.22(10)	N/A(N/A)	23.89(13)	48.46(6)
Ionosphere	73.92(8)	74.88(3)	88.06(3)	72.03(12)	82.91(5)	41.78(14)	41.41(11)	85.87(4)	69.89(10)	65.99(11)	46.58(13)	64.98(12)	73.04(9)	88.61(1)
libras	16.16(11)	30.97(2)	24.69(12)	30.97(2)	24.95(2)	27.47(8)	27.47(8)	27.47(8)	27.47(8)	27.47(8)	24.84(7)	24.84(7)	24.84(7)	24.84(7)
letter	6.86(12)	6.11(14)	34.02(1)	14.85(2)	21.43(4)	8.389(30)	30.2(2)	28.63(3)	6.77(13)	6.94(10)	10.16(7)	11.68(6)	6.87(11)	8.49(8)
Lymphography	97.02(3)	93.59(4)	23.08(11)	97.62(1)	36.68(10)	91.83(5)	38.69(9)	26.65(12)	88.68(7)	90.87(6)	6.53(14)	19.52(13)	44.54(8)	97.31(2)
magic.gam	92.27(6)	51.94(12)	54.76(9)	68.85(4)	54.12(11)	62.41(5)	75.63(1)	67.89(3)	59.18(7)	54.38(10)	35.67(14)	46.92(13)	58.49(8)	64.72(4)
manometry	19.29(25)	12.54(30)	9.76(28)	11.15(10)	11.14(10)	21.31(3)	15.91(6)	13.41(8)	40.67(2)	41.28(1)	6.65(13)	N/A(N/A)	14.75(7)	20.07(6)
manometry	33.2(3)	23.9(1)	28.82(5)	25.51(8)	28.82(5)	11.41(10)	11.41(10)	11.41(10)	11.41(10)	11.41(10)	11.41(10)	11.41(10)	11.41(10)	11.41(10)
musk	99.89(3)	10.61(9)	2.82(13)	10.0(1)	2.61(14)	100(1)	9.65(10)	7.59(11)	34.79(7)	34.95(6)	5.51(2)	32.75(8)	47.65(5)	99.61(4)
optdigits	2.76(14)	2.92(13)	6.06(3)	10.08(1)	4.42(6)	10.03(2)	3.06(12)	4.39(7)	4.36(8)	3.43(10)	3.38(11)	5.59(4)	3.95(9)	5.09(5)
Parkinson	55.71(2)	49.14(6)	39.64(10)	49.65(4)	41.02(9)	33.32(13)	45.39(8)	37.65(12)	37.05(12)	49.3(5)	14.68(14)	53.25(1)	51.29(3)	46.04(7)
Parkinson	71.71(21)	68.64(13)	76.41(10)	82.53(1)	81.44(4)	88.93(1)	78.46(9)	84.28(2)	80.07(6)	68.62(14)	81(5)	69.45(12)	79.74(7)	78.96(8)
pendigits	23.65(3)	23.52(4)	37.88(12)	17.27(8)	2.89(14)	29.27(1)	6.50(1)	4.26(11)	21.22(6)	23.07(5)	2.98(13)	4.67(10)	18.71(7)	26.05(2)
penplot	46.61(4)	50(7)	47.39(12)	44.70(12)	55.66(1)	47.39(12)	50.0(1)	47.39(12)	47.39(12)	47.39(12)	47.39(12)	47.39(12)	47.39(12)	47.39(12)
satellite	59.64(8)	57.61(8)	37.68(13)	61.48(5)	39.71(2)	67.25(1)	63.09(4)	47.23(11)	55.58(9)	65.94(2)	32.38(14)	58.33(7)	61.94(4)	65.92(3)
satimage-2	85.69(3)	82.71(4)	4.29(13)	97.09(1)	8.81(2)	78.04(6)	39.14(9)	26.11(10)	76.55(7)	63.25(8)	2.29(14)	22.07(11)	80.52(5)	93.45(2)
shuttle	92.35(6)	85.29(7)	13.76(12)	60.98(8)	12.17(13)	98.40(3)	30.18(10)	20.17(11)	96.56(2)	95.76(4)	8.28(14)	93.2(5)	48.75(9)	97.62(1)
skin	17.44(11)	19.03(7)	18.26(9)	20.68(1)	16.38(12)	23.7(5)	28.72(1)	28.72(1)	1.98(10)	15.96(13)	20.78(6)	N/A(N/A)	18.44(6)	26.04(3)
skp	66.7(1)	18.9(11)	10.25(9)	26.82(9)	61.13(3)	35.7(2)	66.7(2)	33.7(1)	50.01(5)	55.01(5)	16.7(12)	50.03(4)	16.7(12)	12.81(3)
svm	40.15(5)	35.16(12)	41.18(2)	34.73(1)	50.03(4)	40.03(10)	40.03(10)	40.03(10)	66.8(1)	33.95(2)	55.68(8)	55.68(8)	51.75(5)	51.75(5)
speech	1.97(10)	1.96(11)	2.52(2)	1.99(9)	2.25(4)	2.09(6)	2.02(8)	2.13(5)	1.94(12)	1.77(14)	5.24(1)	2.03(7)	1.79(13)	2.31(3)
Stamps	41.09(3)	31.39(8)	21.12(1)	23.66(9)	16.5(13)	35.24(6)	25.33(10)	20.28(12)	43.1(2)	38.17(5)	10.4(4)	43.72(1)	39.49(4)	39.49(4)
thyroid	44.34(4)	21.23(9)	20.81(10)	25.95(6)	28.5(7)	50.98(3)	34.98(5)	23.56(8)	19.64(11)	54.05(2)	6.85(14)	16.06(12)	14.68(13)	11.61(11)
vertebral	10.49(10)	10.94(8)	14.24(2)	11.58(6)	13.85(4)	9.23(13)	10.57(9)	11.79(8)	8.89(14)	11.24(7)	13.87(15)	15.24(1)	9.68(12)	10.46(11)
vehicle	8.24(12)	8.24(12)	34.72(12)	55.96(2)	12.41(13)	55.96(2)	12.41(13)	55.96(2)	12.41(13)	55.96(2)	12.41(13)	55.96(2)	12.41(13)	12.41(13)
Waveform	5.79(10)	4.37(12)	5.53(14)	18.98(1)	14.11(2)	5.86(9)	13.04(5)	6.66(6)	6.96(6)	6.86(7)	3.42(13)	3.11(14)	4.71(11)	6.24(8)
WBC	82.29(6)	89.87(3)	11.73(2)	92.73(1)	73.73(2)	73.56(6)	66.55(9)	65.54(10)	86.19(4)	15.96(11)	N/A(N/A)	78.77(7)	90.49(2)	90.49(2)
WDBC	75.46(5)	71.88(6)	14.93(14)	79.62(1)	50.52(9)	88.84(1)	43.72(10)	35.6(11)	84.78(2)	57.91(8)	15.37(13)	18.48(12)	66.11(7)	78.53(4)
win	3.13(14)	3.62(12)	5.05(3)	3.64(13)	4.95(8)	4.98(4)	4.73(5)	5.5(11)	3.69(10)	4.14(7)	5.45(2)	4(8)	3.36(13)	4.23(6)
yeast	30.87(1)	23.56(12)	7.77(13)	5.83(14)	8.45(10)	43.08(3)	8.43(11)	7.95(12)	45.71(2)	18.37(7)	8.84(9)	17.51(8)	48.82(1)	25.96(5)
yeast	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)	23.91(5)
year	20.99(1)	29.84(12)	31.64(4)	30.93(7)	31.27(6)	32.75(3)	29.33(14)	29.96(9)	30.71(8)	31.36(5)	33.43(1)	29.92(10)	33.92(2)	29.92(10)
FashionMNIST	31.42(6)	31.97(5)	16.85(13)	38.9(1)	20.73(11)	29.43(8)	33.87(2)	28.72(9)	30.32(7)	32.53(3)	19.72(12)	14.44(4)	27.32(10)	32.35(4)
CIFAR10	10.59(6)	10.19(7)	13.02(1)	10.61(5)	11.61(2)	8.38(13)	11.13(3)	11.06(4)	8.77(12)	9.29(10)	9.8(8)	7.73(14)	9.72(9)	8.97(11)
SVHN	8.66(5)	8.65(6)	9.24(2)	8.87(3)	8.77(3)	7.94(12)	9.46(1)	8.52(8)	7.64(11)	7.82(10)	7.27(14)	7.29(13)	8.74(1)	8.10(9)
Agnews	5.74(8)	5.69(9)	7.02(5)	12.12(7)	6.75(10)	5.88(10)	6.61(3)	4.5(8)	5.09(7)	5.13(9)	5.4(11)	5.7(13)	5.4(11)	5.4(11)
Agnews	8.85(9)	5.64(13)	5.72(11)	6.07(4)	5.97(13)	5.45(12)	6.41(4)	6.08(13)	6.06(5)	5.29(14)	6.58(12)	5.92(8)	5.95(9)	5.95(9)
Imdb	4.55(12)	4.44(14)	4.83(5)	4.75(6)	5.16(1)	4.74(7)	4.49(13)	4.7(9)	4.9(2)	4.9(2)	4.86(4)	4.65(10)	4.59(11)	4.74(7)
Yelp	7.62(12)	7.75(9)	8.52(4)	7.68(10)	8.68(3)	7.81(8)	9.85(1)	9.2(2)	8.01(5)	7.98(6)	5.15(14)	6.72(13)	7.65(11)	7.88(7)

Table 7: AUCROC of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 5\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	58.73(15)	76.43(8)	76.65(7)	81.41(2)	81.14(3)	78.25(6)	75.24(9)	81.11(4)	40.02(16)	79.94(5)	68.28(12)	82.89(1)	67.62(13)	70.21(1)	63.19(14)	73.81(10)
ALOI	56.71(4)	57.88(2)	52.13(9)	47.46(16)	48.51(13)	52.77(8)	65.99(1)	48.37(14)	49.86(11)	48.95(12)	47.91(15)	51.97(10)	56.68(5)	55.28(6)	57.22(3)	55.05(7)
anthyroid	76.22(14)	79.92(11)	78.91(12)	81.64(10)	82.42(7)	89.05(6)	97.24(2)	81.99(9)	49.81(6)	71.31(15)	77.35(13)	52.58(1)	82.55(8)	93.14(5)	95.24(3)	98.51(1)
Arrhythmia	74.58(1)	59.22(9)	70.78(2)	56.94(10)	53.45(13)	60.73(7)	62.84(6)	42.94(16)	47.68(14)	47.48(15)	54.81(2)	56.18(11)	66(1)	65.34(4)	65.19(5)	60.03(8)
breastw	92.28(1)	90.67(12)	97.61(4)	99.18(2)	92.79(9)	92.67(10)	96.72(6)	99.73(1)	72.44(15)	95.42(7)	49.28(16)	88.59(13)	85.59(14)	97.53(5)	95.12(8)	99.1(3)
cardio	82.46(13)	79.71(14)	95.56(5)	97.12(1)	96.36(3)	87.96(10)	95.73(4)	71.04(15)	83.3(12)	90.85(7)	N/A	92.65(6)	85.07(11)	88.03(9)	89.79(8)	96.82(2)
Cardiotocography	54.77(16)	81.24(12)	89.88(3)	91.17(1)	90.73(2)	81.89(10)	83.73(8)	86.77(5)	80.19(13)	84.67(6)	60.64(15)	81.73(11)	79.49(14)	83.03(9)	84.55(7)	89.12(4)
comm.and.crim	63.04(13)	64.47(12)	71.87(4)	73.29(3)	69.83(7)	62.67(14)	70.41(6)	79.86(1)	52.38(15)	68.73(8)	51.26(16)	65.79(11)	68.69(1)	71.2(5)	67.51(10)	76.26(2)
concrete	64.08(4)	60.39(15)	66.53(11)	79.09(2)	78.91(3)	71.96(8)	72.87(7)	70.79(9)	70.34(10)	54.96(16)	65.61(12)	77.79(4)	64.51(13)	77.39(5)	75.75(6)	79.75(1)
cover	44.46(16)	94.43(9)	99.69(2)	99.47(3)	99.47(3)	91.01(11)	97.18(6)	96.92(7)	89.4(12)	96.38(8)	68.67(14)	99.92(1)	84.32(13)	57.48(15)	94.08(10)	99.24(5)
fault	64.26(12)	68.76(4)	68.7(5)	67.54(7)	66.67(9)	62.94(14)	69.3(3)	63.65(13)	62.06(15)	66.07(10)	50.92(16)	66.99(8)	65(11)	69.68(2)	68.15(6)	72.02(1)
glass	67.84(16)	84.34(11)	82.27(8)	91.37(9)	94.46(4)	93.72(5)	95.73(3)	92.73(6)	80.63(14)	82.59(12)	89.96(10)	91.48(8)	78.85(15)	92.17(7)	95.97(2)	98.37(1)
HearDisease	52.15(15)	67.57(13)	74.64(8)	76.82(6)	75.59(7)	67.67(11)	77.94(5)	80.23(3)	48.66(16)	73.44(9)	61.19(14)	68.94(10)	67.65(12)	80.55(2)	78.37(4)	83.57(1)
Hepatitis	56.94(15)	74.83(12)	88.1(4)	86.62(5)	83.71(8)	72.08(13)	85.44(6)	77.79(11)	50(16)	85.11(7)	80.48(10)	82.99(9)	69.71(14)	89.29(3)	89.66(2)	91.72(1)
http	99.81(1)	99.91(0)	99.98(8)	100(1)	100(1)	99.99(7)	98.33(12)	98.33(12)	83.31(14)	0.14(16)	100(1)	100(1)	100(1)	99.87(15)	100(1)	99.97(9)
imgsg	62.24(16)	85.08(4)	79.61(2)	84.44(6)	84.39(7)	83.51(9)	84.53(5)	57.22(16)	87.90(10)	89.03(8)	67.59(15)	92.84(6)	86.55(11)	92.07(8)	94.54(4)	97.32(1)
InternetAds	67.94(8)	72.98(4)	79.79(2)	61.68(11)	57.82(12)	63.57(10)	75.77(1)	51.49(14)	69.28(5)	57.14(13)	50.82(15)	N/A	68.01(7)	68.68(6)	67.87(9)	81.94(1)
ionosphere	92.09(4)	81.01(9)	94.21(1)	68.67(13)	67.54(14)	63.46(15)	93.43(3)	85.02(7)	83.55(8)	71.41(12)	50.49(16)	76.83(10)	74.95(11)	88.77(5)	87.23(6)	93.99(2)
landsat	46.33(16)	85.54(4)	59.84(15)	77.18(2)	78.32(11)	81.22(7)	87.34(2)	71.24(13)	79.27(10)	80.25(8)	66.81(14)	86.98(3)	79.52(9)	84.7(5)	82.67(6)	88.44(1)
letter	69.48(1)	73.09(3)	66.43(5)	65.01(7)	63.18(11)	61.18(8)	83.13(1)	48.83(14)	61.02(16)	42.09(15)	76.96(2)	59.82(13)	69.32(1)	63.92(5)	67.89(7)	83.54(8)
Lymphography	96.84(2)	86.98(3)	92.94(4)	80.25(9)	79.68(10)	93.47(3)	77.41(3)	73.19(14)	50(16)	70.06(15)	87.01(7)	91.87(5)	78.22(12)	91.31(6)	78.66(11)	99.71(1)
magic_gamma	61.44(5)	82.08(6)	78.24(10)	82.76(5)	83.38(4)	84.04(3)	84.58(2)	76.13(13)	49.3(16)	85.08(1)	71.05(14)	81.91(7)	78.06(12)	79.11(9)	78.22(11)	81.85(8)
mammography	75.46(13)	90.96(6)	91.9(3)	92.9(2)	93.17(1)	91.28(4)	87.22(10)	87.26(9)	69.49(16)	89.65(8)	77.88(12)	90.8(7)	72.62(14)	73.12(15)	86.44(11)	91.15(5)
mnist	49.09(14)	84.99(12)	92.3(7)	93.5(3)	93.3(5)	86.69(5)	91.28(1)	57.22(16)	87.90(10)	89.03(8)	67.59(15)	92.84(6)	86.55(11)	92.07(8)	94.54(4)	97.32(1)
musk	73.94(9)	96.46(10)	93.97(11)	100(1)	100(1)	100(1)	99.95(3)	83.33(13)	54.07(16)	100(1)	78.75(15)	99.89(6)	81.61(14)	83.97(12)	99.66(8)	99.78(7)
optdigits	46.44(16)	94.98(10)	99.67(4)	99.98(1)	99.98(1)	99.94(3)	98.23(8)	80.46(13)	90.94(12)	99.48(5)	52.71(15)	98.84(7)	93.91(11)	78.07(14)	96.75(9)	99.04(6)
PageBlocks	72.69(15)	93.01(4)	91.22(6)	86.42(12)	90.38(7)	89.67(9)	94.73(2)	88.46(10)	58.16(16)	90.15(8)	82.94(13)	88.03(11)	76.84(14)	92.74(5)	94.23(3)	96.8(1)
Parkinson	50.21(16)	75.56(9)	69.54(12)	83.63(8)	80.94(6)	74.14(11)	81.76(5)	87.09(1)	53.05(15)	74.92(10)	63.65(14)	80.55(7)	65.13(13)	82.1(4)	79.14(8)	86.56(2)
pendigits	56.82(16)	97.41(11)	99.75(2)	99.76(1)	99.69(3)	99.59(4)	98.28(8)	99.1(5)	84.95(12)	98.23(10)	70.59(14)	98.27(9)	84.29(13)	64.95(15)	98.33(7)	99.03(6)
Pima	60.11(3)	63.23(12)	73.18(4)	78.02(1)	76.43(2)	63.55(11)	69.92(8)	71.13(6)	47.99(15)	32.37(16)	51.59(14)	68.69(9)	67.36(10)	71.8(5)	70.8(7)	73.8(3)
satellite	69.49(15)	88.65(3)	81.31(12)	84.34(10)	84.41(9)	85.32(7)	89.34(2)	79.37(13)	61.94(16)	87.79(5)	75.45(14)	87.82(4)	82.71(11)	80.74(6)	85.28(8)	90.22(1)
satimage-2	96.81(0)	98.11(3)	97.79(7)	97.82(6)	97.82(6)	98.55(2)	97.83(5)	91.43(12)	65.61(15)	96.65(9)	85.08(14)	96.99(8)	87.97(13)	90.02(16)	93.66(11)	98.33(1)
shuttle	77.41(16)	98.71(4)	98.83(3)	97.57(8)	97.09(7)	97.29(11)	77.73(1)	94.63(14)	79.84(13)	97.49(9)	97.97(6)	97.44(10)	78.01(15)	97.16(12)	98(5)	98.85(2)
skin	52.78(16)	99.41(4)	90.56(14)	95.72(11)	95.24(12)	98.56(8)	99.22(5)	54.09(13)	99.6(3)	97.23(10)	99.83(1)	99.06(6)	78.23(15)	99(7)	98.53(9)	99.72(2)
smtp	51.69(14)	82.74(3)	76.47(8)	75.52(10)	74.57(11)	79.75(10)	89.78(5)	80.81(9)	68.44(14)	87.16(6)	54.44(15)	87.75(5)	78.45(12)	84.84(7)	82.86(10)	91.57(1)
Spambase	53.66(16)	70.72(13)	83.55(8)	90.24(3)	90.87(2)	79.04(11)	89.24(4)	79.84(10)	68.44(14)	87.16(6)	54.44(15)	87.75(5)	78.45(12)	84.84(7)	82.86(10)	91.57(1)
speech	47.49(15)	53.13(10)	62.16(13)	75.91(7)	65.61(17)	85.24(2)	49.88(14)	49.88(14)	54.66(16)	75.91(11)	71.39(15)	91.02(6)	73.95(13)	84.9(8)	87.03(7)	96.07(1)
Stamps	72.13(4)	74.61(12)	95.3(3)	94.14(5)	96.07(1)	79.4(10)	89.04(4)	83.22(9)	54.66(16)	75.91(11)	71.39(15)	91.02(6)	73.95(13)	84.9(8)	87.03(7)	96.07(1)
thyroid	92.65(1)	95.21(10)	99.55(1)	99.5(2)	99.5(2)	99.42(4)	98.98(7)	96.53(8)	71.37(14)	85.18(12)	51.29(16)	99.35(5)	78.61(13)	71.21(5)	95.76(9)	98.99(6)
vertebral	37.14(16)	61.49(13)	55.86(15)	77.65(4)	77.24(5)	77.2(7)	69.29(11)	81.21(3)	71.9(9)	67.05(12)	69.77(10)	81.88(1)	56.29(14)	73.5(8)	81.65(2)	77.22(6)
vowels	78.54(9)	80.31(7)	86.47(6)	89.13(2)	89.13(2)	86.69(5)	91.82(1)	73.29(12)	55.81(13)	14.35(16)	80.28(14)	82.81(1)	66.06(14)	73.15(13)	73.74(13)	87.13(4)
Waveform	52.99(16)	68.05(10)	80.58(5)	84.84(1)	82.3(3)	64.09(13)	74.68(7)	77.42(6)	60.62(15)	61.03(14)	72.6(8)	80.85(4)	64.74(12)	64.81(11)	71.76(9)	84.82(2)
WBC	92.78(7)	94.96(6)	97.88(5)	98.38(2)	98.28(3)	68.09(14)	74.32(12)	65(15)	58.91(16)	76.82(10)	72.86(13)	97.98(4)	81.68(8)	81.23(9)	76.14(11)	98.61(1)
WDBC	97.52(7)	90.44(11)	98.81(5)	99.91(1)	99.82(3)	99.73(4)	95.24(10)	56.02(15)	61.79(14)	0.45(16)	82.49(13)	99.87(2)	90.15(12)	95.27(9)	95.46(8)	98.49(6)
Wilt	83.88(16)	81.89(4)	49.93(15)	67.91(12)	68.57(10)	83.57(8)	85.71(6)	74.09(10)	65.91(13)	78.45(3)	59.85(16)	82.81(1)	61.11(16)	66.07(15)	75.52(10)	89.99(2)
wine	65.76(4)	82.24(12)	99.91(3)	100(1)	100(1)	99.82(4)	98.29(11)	95.97(2)	50(15)	3.91(16)	88.46(10)	99.82(4)	73.73(13)	89.63(9)	94.27(8)	99.2(6)
WPBC	47.66(15)	55.28(13)	57.32(10)	68.74(2)	69.92(1)	57(11)	61.95(6)	55.92(12)	44.53(16)	54.06(14)	64.63(3)	60.63(7)	60.11(8)	57.48(9)	62.26(5)	64.39(4)
yeast	48.22(15)	54.78(12)	47.57(16)	65.99(2)	66.12(1)	59.7(7)	55.83(10)	63.75(3)	62.62(4)	53.76(13)	50.93(14)	60.65(6)	56.28(9)	57.22(8)	55.18(11)	62.04(5)
FashionMNIST	79.87(16)	86.83(6)	87.79(67)	88.56(1)	87.79(67)	88.56(1)	87.79(67)	88.56(1)	87.79(67)	88.56(1)	87.79(67)	88.56(1)	87.79(67)	88.56(1)	87.79(67)	88.56(1)
CIFAR10	60.37(13)	64.91(9)	71.88(4)	71.06(5)	65.98(8)	67.83(7)	75.5(2)	50.54(16)	61.36(11)	64.51(10)	55.79(14)	61.16(12)	54.31(15)	70.83(6)	75.67(1)	73.73(3)
SVHN	57.94(12)	59.56(10)	64.82(4)	67.36(1)	62.21(9)	63.69(6)	67.23(2)	49.84(16)	62.54(8)	58.4(11)	57.44(13)	56.13(14)	52.59(15)	63.56(7)	66.56(3)	69.93(5)
Agnews	59.61(14)	71.91(11)	84.3(3)	84.67(2)	84.73(1)	76(7)	78.14(5)	52.97(15)	75.65(8)	82.97(4)	64.24(13)	67.08(12)	52.85(16)	73.49(10)	78.08(6)	75.48(9)
Amazon	58.78(14)	61.38(12)	69.54(8)	72.33(6)	72.33(6)	74.7(5)	71.41(5)	72.17(1)	74.51(13)	78.45(3)	59.85(16)	61.11(16)	66.07(15)	75.52(10)	76.52(9)	80.99(2)
Imdb	49.59(16)	59.52(13)	72.56(5)	77.09(2)	72.52(6)	69.03(9)	72.52(6)	72.64(14)	72.84(13)	73.43(3)	65.74(10)	60.16(12)	51.58(15)	63.62(11)	69.74(8)	71.08(7)
Yelp	67.87(12)	64.79(13)	79.86(7)	86.8(1)	86.73(2)	82.74(4)	80.12(6)	58.21(15)	76.84(8)	86.15(3)	63.24(14)	70.17(11)	53.62(16)	75.66(9)	80.88(5)	74.69(10)

Table 8: AUCPR of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 5\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	59.33(15)	72.49(8)	77.69(6)	80.71(2)	80.49(3)	75.04(7)	72.02(9)	79.87(4)	46.51(16)	79.61(5)	66.38(11)	82.21(1)	64.22(13)	66.13(12)	63.05(14)	70.24(10)
ALOI	53.85(1)	4.359(7)	4.359(7)	4.55(8)	4.55(8)	4.359(7)	7.28(1)	2.49(16)	3.29(16)	3.29(16)	3.29(16)	4.77(16)	6.2(16)	4.77(16)	6.2(16)	5.49(16)

Table 13: AUROC of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 50\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	63.41(16)	84.97(3)	79.23(14)	81.95(8)	81.88(9)	83.54(6)	83.95(5)	81.23(11)	71.16(15)	85.64(2)	80.21(13)	86.19(1)	82.17(7)	81.3(10)	80.36(12)	84.25(4)
ALOI	54.52(11)	64.14(6)	53.63(12)	52.88(14)	52.95(13)	57.95(7)	79.27(1)	52.7(15)	52.07(16)	57.46(8)	56.16(9)	54.72(10)	78.9(2)	70.12(5)	70.28(3)	70.28(3)
anthyroid	76.85(16)	95.54(9)	83.31(11)	82.81(13)	82.58(14)	91.43(10)	99.21(5)	83.01(12)	82.03(15)	96.62(8)	98.13(7)	98.97(6)	99.34(3)	99.38(1)	99.37(2)	99.27(4)
Arrhythmia	75.49(15)	85.49(9)	88.53(6)	82.94(10)	82.54(11)	82.34(12)	92.04(5)	55.5(16)	70.77(14)	81.49(13)	86.77(8)	87.88(7)	93.23(4)	94.08(1)	93.38(2)	93.36(2)
breastw	94.15(16)	98.44(13)	98.84(12)	99.67(2)	99.59(4)	98.97(11)	99.17(9)	99.75(1)	98.32(14)	99.62(3)	95.52(15)	99.1(10)	99.19(8)	99.42(5)	99.46(6)	99.38(7)
cardio	85.52(15)	96.64(12)	98.68(6)	98.75(4)	98.62(7)	96.56(13)	99.38(1)	94.69(14)	97.54(10)	96.88(11)	NANA	98.7(5)	99.31(2)	98.52(8)	98.28(9)	99.18(3)
Cardiotocography	58.27(16)	94.69(11)	95.86(7)	95(10)	95.22(9)	94.3(12)	96(4)	93.07(14)	94.29(13)	95.79(7)	89.09(15)	95.52(8)	96.23(3)	96.42(2)	95.95(5)	96.53(1)
comm.and.crim	68.29(16)	81.81(13)	87.42(7)	91.09(1)	90.82(2)	79.85(14)	88.72(5)	87.66(1)	83.59(12)	87.28(8)	78.43(15)	83.79(11)	89.79(4)	87.25(9)	86.95(10)	90.54(3)
concrete	64.72(16)	84.84(12)	75.95(15)	88.62(10)	89.21(7)	88.74(9)	90.15(6)	86.27(14)	83.7(13)	92.56(2)	87.07(11)	90.18(5)	88.78(8)	90.67(3)	90.21(4)	92.62(1)
cover	42.11(16)	99.95(1)	99.95(1)	99.93(7)	99.94(4)	99.88(8)	99.94(3)	99.81(12)	98.99(14)	99.86(9)	97.07(15)	99.94(4)	99.95(1)	99.83(10)	99.82(11)	99.94(4)
fault	66.32(16)	79.36(7)	76.32(11)	75.61(12)	77.47(9)	75.53(13)	81.29(4)	69.99(15)	76.87(10)	80.56(5)	74.97(14)	78.64(8)	83.03(2)	81.77(3)	80.15(6)	84.08(1)
glass	68.48(16)	95.63(10)	90.11(13)	89(14)	90.35(12)	99.61(6)	99.66(5)	94.38(11)	98.52(9)	85.43(15)	99.83(2)	99.67(4)	99.73(3)	99.34(8)	99.41(7)	99.87(1)
HearDisease	61.49(16)	92.44(11)	90.27(14)	92.81(10)	92.84(9)	90.96(13)	97.25(4)	92.09(12)	89.49(15)	94.48(7)	93.48(8)	96.92(5)	96.84(6)	97.84(1)	97.71(3)	97.74(2)
Hepatitis	66.13(16)	99.67(11)	99.74(10)	99.23(12)	98.83(13)	99.96(8)	99.94(9)	97.18(14)	94.81(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
http	99.77(11)	100(1)	99.98(8)	100(1)	100(1)	99.93(10)	98.33(12)	98.33(12)	83.33(14)	0.39(16)	100(1)	100(1)	100(1)	40.96(15)	100(1)	99.97(9)
imgsg	71.32(16)	97.29(8)	84.71(4)	87.58(13)	87.85(12)	92.45(10)	97.53(7)	80.55(15)	92.41(11)	94.45(9)	97.64(6)	97.81(5)	98.53(2)	98.25(3)	98.05(4)	99.21(1)
InternetAds	68.93(15)	90.13(9)	94.33(2)	92.95(5)	92.47(7)	89.31(12)	92.35(8)	76.24(14)	89.42(11)	89.53(10)	80.04(13)	NANA	95.73(3)	92.93(6)	93.48(4)	95.37(1)
Ionosphere	92.84(11)	96.83(7)	97.55(6)	91.77(14)	92.28(12)	92.26(13)	98.93(2)	89.73(15)	94.26(10)	96.68(8)	89.12(16)	96.17(9)	97.85(5)	99.04(1)	98.83(3)	98.77(4)
landsat	55.06(16)	94.52(4)	60.57(15)	80.19(12)	79.86(13)	87.52(11)	94.6(3)	71.73(14)	90.24(10)	92.93(8)	91.23(9)	93.96(7)	95.09(2)	94.47(5)	94.45(6)	95.46(1)
letter	70.16(16)	83.75(10)	88.74(6)	84.17(9)	83.32(11)	82.56(13)	91.91(2)	76.71(15)	81.81(4)	82.59(12)	84.9(8)	89.13(5)	90.57(4)	91.58(3)	87.49(7)	93.25(1)
Lymphography	98.18(14)	99.88(13)	100(1)	100(1)	100(1)	100(1)	100(1)	99.92(12)	96.82(15)	100(1)	93.94(16)	100(1)	100(1)	100(1)	100(1)	100(1)
magic.gamma	56.32(16)	90.24(4)	82.46(14)	82.94(13)	83.19(12)	84.89(10)	90.13(6)	74.73(15)	83.81(11)	90.41(3)	85.29(9)	88.45(7)	90.68(2)	90.18(5)	88.42(8)	91.07(1)
mammography	77.28(15)	94.42(2)	93.17(11)	93.28(9)	93.29(8)	93.93(5)	94.12(3)	91.32(12)	70.71(6)	93.28(9)	93.62(6)	93.92(4)	93.47(7)	87.29(14)	91.21(13)	94.48(1)
mnist	99.14(7)	99.36(4)	99.01(9)	99.18(8)	99.13(9)	100(1)	100(1)	99.02(15)	99.01(12)	98.42(11)	95.66(14)	99.66(2)	99.56(1)	99.31(5)	99.54(1)	99.51(2)
mnus	99.90(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	99.43(14)	99.41(15)	100(1)	99.98(13)	100(1)	100(1)	100(1)	100(1)	100(1)
optdigits	46.9(16)	99.97(6)	99.99(1)	99.99(1)	99.99(1)	99.99(1)	99.93(7)	98.05(14)	99.8(13)	99.98(5)	97.59(15)	99.91(10)	99.91(10)	99.91(10)	99.92(8)	99.92(8)
PageBlocks	78.46(16)	97.21(8)	94.08(12)	88.09(15)	90.53(14)	96.22(9)	98.31(4)	91.83(13)	94.17(11)	95.63(9)	97.49(7)	97.61(6)	98.29(5)	98.39(2)	98.34(3)	98.87(1)
Parkinson	57.47(16)	97.28(8)	79.15(15)	90.16(10)	88.32(12)	83.85(14)	98.64(6)	88.05(13)	89.72(11)	91.89(9)	99.01(5)	99.23(1)	97.94(7)	99.25(3)	99.24(4)	99.31(2)
pendigits	59.36(16)	99.99(2)	99.88(3)	99.76(13)	99.81(11)	99.87(9)	99.93(6)	99.69(14)	99.77(12)	99.83(10)	98.71(5)	99.97(4)	99.99(2)	99.95(5)	99.91(7)	100(1)
Pima	63.66(16)	80.63(11)	75.75(15)	82.93(9)	82.73(10)	80.5(12)	88.62(5)	83.54(7)	77.38(14)	82.97(8)	80.45(13)	84.25(6)	89.73(1)	89.63(2)	88.73(4)	89.13(3)
satellite	75.12(16)	95.28(7)	81.26(14)	85.38(12)	84.85(13)	86.32(11)	95.55(3)	78.88(15)	91.48(8)	90.66(10)	91.13(9)	95.39(6)	96.25(1)	95.42(4)	95.45(5)	96.19(2)
satellite-2	88.33(13)	99.28(3)	99.45(1)	99.06(4)	99.05(5)	98.31(7)	98.5(11)	98.10(1)	97.38(12)	98.26(8)	97.43(11)	98.16(9)	96.78(14)	95.09(16)	96.16(15)	98.59(6)
shuttle	65.45(16)	98.97(6)	98.91(7)	97.57(14)	97.67(12)	98.42(9)	99.98(1)	97.37(15)	92.88(10)	97.63(13)	98.86(8)	98.28(1)	99.32(5)	99.85(3)	99.65(6)	99.62(2)
skin	52.08(16)	99.92(3)	89.02(15)	95.7(12)	95.35(13)	98.77(11)	99.88(6)	93.94(14)	99.73(10)	99.88(6)	99.93(2)	99.9(4)	99.77(9)	99.89(5)	99.87(8)	99.94(1)
smtp	56.11(14)	99.22(4)	92.43(6)	87.46(7)	84.95(10)	84.24(11)	98.52(5)	50.48(15)	66.96(13)	86.42(9)	83.31(2)	87.3(8)	100(1)	43.94(16)	99.72(3)	99.99(2)
Spambase	57.07(16)	94.37(8)	91.77(14)	92.05(12)	94.09(9)	93.64(11)	97.38(2)	83.88(15)	91.82(13)	97.05(4)	93.71(10)	95.27(3)	95.74(6)	97.15(3)	96.55(5)	97.69(1)
speech	59.94(14)	99.94(14)	73.34(7)	81.56(14)	80.13(15)	79.89(6)	67.42(1)	70.74(16)	95.56(8)	66.84(14)	72.88(15)	90.89(6)	96.63(5)	66.99(16)	70.35(9)	86.84(3)
Stamps	74.37(16)	98.03(14)	98.84(11)	98.41(13)	99(9)	98.85(10)	99.61(4)	99.07(8)	97.3(15)	99.12(7)	98.47(12)	99.65(3)	99.32(6)	99.67(2)	99.56(5)	99.72(1)
thyroid	92.97(16)	99.04(13)	99.76(5)	99.76(5)	99.76(5)	99.76(5)	99.9(1)	99.62(10)	95.95(15)	99.57(11)	97.46(14)	99.7(9)	99.22(12)	99.8(4)	99.85(3)	99.86(2)
vertebral	39.88(16)	88.41(0)	79.82(14)	81.09(13)	82.54(12)	92.65(8)	98.71(1)	87.03(11)	91.13(9)	74.28(15)	96.42(16)	96.29(7)	97.79(5)	97.87(4)	98.22(3)	98.59(2)
vowels	92.21(16)	92.21(16)	98.08(9)	98.91(4)	98.91(4)	98.91(4)	98.13(7)	98.13(7)	92.21(16)	98.13(7)	92.21(16)	92.21(16)	92.21(16)	98.85(12)	97.67(1)	98.85(12)
Waveform	53.02(16)	85.4(14)	91.74(4)	91.02(7)	91.47(5)	86.57(13)	91.05(6)	86.98(12)	91.97(2)	94.73(1)	75.22(15)	90.3(9)	87.11(11)	90.44(8)	89.69(10)	91.77(3)
WBC	95.21(13)	98.38(10)	99.22(6)	99.43(4)	99.14(7)	92.44(15)	99.9(1)	98.95(9)	97(11)	96.5(12)	91.37(16)	99.33(5)	99.68(3)	93.46(14)	99.06(8)	99.83(2)
WDBC	97.76(15)	99.72(12)	100(1)	100(1)	100(1)	100(1)	99.29(13)	100(1)	98.75(14)	100(1)	91.31(16)	100(1)	99.99(9)	99.8(11)	99.89(10)	100(1)
Wilt	62.78(16)	98.38(2)	61.62(15)	68.09(13)	68.09(13)	90.1(10)	98.68(8)	82.18(11)	95.96(8)	66.84(14)	98.62(4)	90.9(9)	99.67(5)	93.03(16)	93.03(16)	93.03(16)
wine	70.5(16)	99.89(14)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	98.73(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
WPBC	49.03(16)	89.94(8)	77.94(14)	82.2(12)	82.5(11)	85.95(10)	97.57(1)	70.04(15)	81.13(3)	86.21(9)	93.25(6)	94.19(5)	92.37(7)	94.79(4)	95.72(2)	95.44(3)
yeast	68.43(16)	67.62(12)	45.2(16)	68.43(16)	68.58(19)	71.08(4)	65.6(13)	65.68(13)	55.86(14)	69.74(8)	69.76(7)	74.27(2)	70.73(5)	71.11(3)	70.67(6)	75.62(1)
FashionMNIST	91.68(7)	94.68(7)	96.32(7)	96.32(7)	96.32(7)	96.32(7)	96.16(16)	96.16(16)	96.16(16)	96.16(16)	96.16(16)	96.16(16)	96.16(16)	96.16(16)	96.16(16)	96.16(16)
CIFAR10	61.15(15)	77.08(11)	84.28(5)	83.35(1)	82.31(7)	79.64(9)	85.06(4)	60.77(16)	80.99(8)	79.19(10)	66.86(14)	71.13(12)	69.51(13)	85.17(3)	84.27(6)	83.53(2)
SVHN	57.85(15)	74.64(10)	80.99(2)	82.12(1)	80.34(3)	77.3(9)	79.31(6)	51.66(16)	77.99(8)	73.79(11)	69.37(12)	68.41(13)	65.37(14)	79.61(5)	78.57(7)	83.33(2)
Agnews	59.67(16)	88.54(12)	92.29(3)	92.15(5)	92.19(4)	90.06(10)	92(8)	66.21(5)	92.11(6)	95.04(1)	89.53(11)	87.88(13)	70.95(14)	91.90(9)	92.03(7)	92.6(2)
Amazon	57.41(15)	78.71(11)	87.47(3)	80.13(12)	80.13(12)	86.94(7)	91.06(13)	81.14(10)	90.6(11)	87.23(14)	90.6(11)	77.21(13)	86.21(14)	86.40(15)	86.40(15)	86.40(15)
Imdb	49.26(16)	80.25(12)	81.63(10)	84(2)	83.72(4)	82.06(8)	82.73(7)	59.97(14)	82.86(13)	81.33(11)	72.21(13)	59.71(5)	91.93(9)	83.19(5)	83.93(3)	83.93(3)
Yelp	68.41(14)	83.38(13)	92.12(2)	92.1(3)	90.53(5)	88.92(10)	90.19(7)	59.38(16)	91.37(4)	94.62(1)	86.64(11)	85.83(12)	66.45(15)	87.91(9)	89.95(8)	90.23(6)

Table 14: AUCPR of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 50\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTstranformer	RF	LGB	XGB	CatB
abalone	63.96(16)	81.08(8)	79.33(10)	81.35(5)	81.18(7)	81.3(6)	81.59(4)	80.16(9)	75.41(15)	85.26(1)	76.16(13)	84.78(2)	78.46(11)	77.85(12)	75.83(14)	81.86(3)
ALOI	3.8(16)	7.79(6)	4.68(11)	4.29(13)	4.29(13)	5.81(6)	14.4(23)	4.04(12)	5.48(12)	5.29(10)	4.43(12)	4.29(13)	12.29(4)	12.26(5)	12.26(5)	15.31(16)
anthyroid	66.48(16)	69.99(9)	43.17(15)	45.48(12)	45.07(14)	56.75(11)	85.86(5)	45.50(15)	58.10(1)	71.27(8)	77.57(7)	83.77(6)	89.09(3)	89.71(1)	89.62(2)	86.92(4)
breast	77.49(16)	87.99(2)	85.32(10)	85.32(10)	85.32(10)	85.32(10)	85.32(10)	46.33(16)	85.32(10)	85.32(10)	85.32(10)	85.32(10)	85.32(10)	85.32(10)	85.32(10)	85.32(10)
breastw	88.43(16)	96.18(13)	97.24(12)	99.33(2)	99.12(4)	98.438(8)	97.7(1)	95.92(1)	96.17(14)	99.2(3)	94.77(15)	97.73(10)	98.08(9)	89.61(6)	95.88(7)	98.67(5)
cardio	44.1(5)	86.44(12)	93.47(2)	91.4(7)	91.59(6)	87.71(13)	93.69(1)	79.41(14)	82.9(13)	87.82(10)	NA/NA	91.72(5)	93.46(3)	87.83(9)	87.77(8)	93.07(4)
Cardiography	37.05(16)	84.1(11)	87.98(2)	84.13(10)	84.69(8)	80.73(13)	86.08(6)	78.74(14)	84.36(9)	87.26(4)	76.85(15)	82.05(12)	87.86(3)	87.24(5)	85.69(7)	88.35(1)
concrete	64.42(15)	89.42(5)	85.1(13)	86.2(13)	85.1(13)	87.7(10)	87.1(10)	87.7(10)	87.7(10)	87.7(10)	87.7(10)	87.7(10)	87.7(10)	87.7(10)	87.7(10)	87.7(10)
cover	65.43(16)	81.13(14)	77.62(15)	86.2(9)	87.92(8)	88.7(10)	88.7(5)	81.36(13)	82.09(12)	91.5(9)	85.63(10)	85.42(11)	88.02(7)	89.21(3)	88.82(4)	91.21(2)
glass	0.87(16)	96.9(1)	96.33(3)	96.42(2)	95.84(4)	94.01(6)	87.82(12)	89.77(10)	81.19(14)	91.629(8)	83.88(13)	92.57(7)	94.82(5)	80.36(15)	89.69(11)	91.8(8)
landsat	52.78(16)	64.98(8)	63.11(13)	61.39(14)	65.97(7)	63.58(12)	69.73(3)	53.82(15)	63.97(10)	74.58(5)	63.79(14)	64.539(70)	70.6(2)	69.36(4)	66.37(6)	73.33(1)
plant	46.49(13)	89.61(13)	86.2(13)	86.2(13)	86.2(13)	86.2(13)	86.2(13)	44.98(16)	86.2(13)	86.2(13)	86.2(13)	86.2(13)	86.2(13)	86.2(13)	86.2(13)	86.2(13)
HeartDisease	54.69(16)	89.12(8)	88.76(15)	91.52(10)	91.61(9)	89.23(13)	96.97(4)	90.2(15)	88.95(14)	93.68(7)	92.69(8)	96.46(6)	95.55(1)	97.7(7)	97.33(3)	97.47(2)
Hepatitis	31.56(16)	98.42(11)	99.10(10)	93.79(12)	91.56(13)	99.79(8)	99.69(9)	86.73(14)	83.17(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
http	67.59(13)	100(1)	99.99(11)	100(1)	100(1)	100(1)	96.68(8)	96.88(14)	67.78(14)	0.37(16)	100(1)	99.44(7)	97.92(7)	93.85(15)	99.41(10)	94.12(10)
Intersect	96.61(16)	85.97(12)	87.92(13)	87.92(13)	87.92(13)	87.92(13)	87.92(13)	92.75(14)	85.97(12)	85.97(12)	85.97(12)	85.97(12)	85.97(12)	85.97(12)	85.97(12)	85.97(12)
Ionosphere	49.41(14)	80.95(9)	87.79(2)	86.41(3)	86.22(4)	74.85(11)	85.8(1)	42.15(15)	83.85(7)	74.61(2)	69.61(13)	NA/NA	78.89(10)	83.73(8)	85.87(5)	91.04(1)
Internad	90.45(15)	95.8(8)	96.72(6)	91.65(13)	91.95(12)	92.4(11)	88.5(6)	85.54(16)	92.78(10)	95.84(7)	90.52(14)	94.829(9)	97.32(5)	98.52(5)	98.39(3)	98.14(3)
landsat	23.29(16)	83.56(5)	41.88(14)	56.1(12)	55.34(13)	62.72(11)	84.47(3)	29.7(15)	77.68(9)	79.45(8)	75.54(10)	80.63(7)	86.35(2)	83.95(4)	83.06(6)	86.69(1)
lets	32.56(16)	64.34(15)	28.06(14)	31.2(13)	31.2(13)	31.2(13)	31.2(13)	14.4(16)	31.2(13)	31.2(13)	31.2(13)	31.2(13)	31.2(13)	31.2(13)	31.2(13)	31.2(13)
Lymphography	74.88(16)	97.82(15)	97.82(15)	100(1)	100(1)	100(1)	100(1)	98.61(13)	76.18(15)	90.2(1)	94.16(14)	100(1)	100(1)	100(1)	100(1)	100(1)
magic-gamma	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.05(5)	76.64(10)	72.86(14)	74.75(13)	74.91(2)	85.48(4)	69.93(15)	79.92(9)	86.86(2)	75.9(11)	80.67(8)	87.57(3)	84.85(6)	81.78(7)	86.87(1)
mammography	42.51(16)	64.0														

Table 15: AUROC of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 75\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	67.09(16)	86.35(2)	77.67(15)	82.15(12)	82.16(11)	81.78(13)	85.15(6)	81.29(14)	85.93(5)	85.99(4)	83.72(9)	86.51(1)	84.5(7)	83.96(8)	82.93(10)	86.11(3)
ALOI	56.22(11)	66.49(6)	51.94(15)	52.85(14)	50.82(16)	63.28(7)	80.73(1)	53.69(12)	58.62(9)	59.71(8)	57.35(10)	53.21(13)	79.83(2)	74.51(3)	73.38(5)	73.46(4)
anthyroid	81.51(16)	97.38(10)	83.51(13)	82.15(15)	83.14(14)	93.21(11)	99.38(5)	84.48(12)	98.1(8)	97.96(9)	98.61(7)	99.41(6)	99.46(2)	99.44(3)	99.47(1)	99.42(4)
Arrhythmia	75.77(16)	92.07(7)	90.85(5)	86.45(13)	86.62(12)	85.85(14)	96.71(3)	82.38(15)	87.01(11)	89.41(10)	91.85(8)	93.72(6)	96.49(4)	96.78(2)	96.46(5)	96.82(1)
breastw	95.04(6)	99.22(14)	99.46(12)	99.73(5)	99.74(4)	98.76(15)	99.68(10)	99.69(9)	99.55(11)	99.75(3)	99.39(13)	99.7(7)	99.79(1)	99.7(7)	99.71(6)	99.78(2)
cardio	87.07(15)	98.39(11)	99.27(7)	98.8(9)	98.56(10)	97.84(12)	99.57(1)	94.77(14)	99.25(8)	97.39(13)	N/A	99.42(4)	99.48(2)	99.35(5)	99.34(6)	99.47(3)
Cardiotocography	60.68(16)	86.59(8)	96.11(10)	95.11(13)	95.25(12)	94.09(15)	97.64(4)	94.31(14)	96.57(9)	96.93(7)	95.84(11)	97.48(5)	97.96(2)	97.79(3)	97.45(6)	98.13(1)
comm.and.crimine	72.22(16)	86.16(13)	88.73(10)	91.39(2)	91.09(4)	83.67(14)	90.06(7)	88.22(11)	91.35(3)	90.32(6)	83.46(15)	86.52(12)	90.84(5)	89.85(8)	89(9)	91.48(1)
concrete	66.39(16)	88.05(13)	78.72(15)	89.6(11)	89.82(10)	88.18(12)	93.54(4)	83.67(14)	92.43(9)	93.34(7)	92.54(8)	93.84(2)	93.74(5)	93.69(3)	93.46(6)	95.63(1)
cover	42.57(16)	99.98(1)	99.95(7)	99.95(7)	99.96(4)	99.8(13)	99.72(15)	99.89(12)	99.97(2)	99.95(7)	99.96(4)	99.94(11)	99.96(4)	99.95(7)	99.78(14)	99.97(2)
fault	67.43(16)	81.36(8)	77.04(12)	75.54(14)	78.07(11)	76.71(13)	83.15(6)	69.77(15)	80.51(9)	82.96(6)	78.39(10)	82.07(7)	86.7(2)	85.46(3)	84.68(4)	86.93(1)
glass	65.51(16)	97.73(10)	88.86(14)	88.61(5)	90.29(3)	99.65(8)	99.98(5)	93.74(11)	99.14(9)	90.53(12)	99.96(7)	99.99(4)	100(1)	100(1)	99.97(6)	100(1)
HearDisese	69.22(16)	95.4(8)	92.37(14)	92.72(13)	93.18(11)	91.98(15)	98.39(5)	93.16(12)	94.04(10)	94.92(9)	97.8(6)	97.55(7)	99.19(1)	98.77(3)	98.74(4)	98.88(2)
Hepatitis	69.27(16)	99.92(9)	99.78(11)	98.86(14)	99.36(12)	99.14(13)	99.93(8)	97.1(15)	99.85(10)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
http	99.81(2)	100(1)	99.98(9)	100(1)	100(1)	99.93(11)	98.33(13)	98.33(13)	100(1)	66.82(16)	100(1)	100(1)	100(1)	78.47(15)	100(1)	99.97(10)
imgseg	76.93(16)	98.14(8)	86.69(14)	88.06(13)	88.08(12)	93.72(11)	98.8(6)	80.96(15)	95.28(10)	95.47(9)	98.8(6)	98.96(5)	99.45(2)	99.24(3)	99.15(4)	99.65(1)
InternetAds	96.23(15)	93.9(11)	96.22(2)	96.23(1)	95.92(4)	93.34(12)	94.86(7)	91.4(16)	97.8(10)	98.21(9)	97.36(12)	98.78(6)	99.15(4)	98.99(5)	99.21(2)	99.21(2)
Ionosphere	93.59(15)	98.34(7)	98.27(8)	94.44(13)	94.08(14)	97.6(11)	99.24(1)	71.73(14)	93.07(10)	94.11(8)	93.36(9)	94.86(7)	96.03(3)	96.16(2)	95.89(4)	96.24(1)
landsat	57.21(16)	95.27(6)	63.11(15)	80.53(12)	80.14(13)	86.7(11)	95.39(5)	75.45(15)	85.08(13)	87.6(9)	91.66(6)	90.85(7)	95.39(2)	93.39(4)	93.86(3)	94.90(1)
letter	70.42(16)	86.83(11)	88.48(6)	87.56(10)	86.46(12)	84.29(14)	93.35(5)	94.97(15)	99.56(4)	99.32(8)	97.75(13)	99.04(11)	99.73(1)	99.55(5)	99.56(2)	99.62(3)
Lymphography	98.31(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	99.71(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
magic.gamma	58.56(16)	90.86(6)	83.36(12)	83.14(14)	83.19(13)	84.17(11)	90.97(4)	77.08(15)	90.47(8)	90.88(5)	88.91(10)	89.92(9)	91.88(2)	91.64(3)	90.7(6)	92.25(1)
mammography	75.53(16)	95.71(1)	93.01(9)	93.1(8)	92.81(12)	95.62(4)	95.71(1)	92.82(11)	81.77(15)	93.34(7)	92.92(10)	95.22(5)	93.46(6)	90.41(14)	92.67(13)	95.67(3)
mnist	99.48(7)	99.31(10)	99.31(10)	99.32(8)	99.32(8)	99.53(6)	99.53(6)	94.97(15)	99.56(4)	99.32(8)	97.75(13)	99.04(11)	99.73(1)	99.55(5)	99.56(2)	99.62(3)
mnist	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
optdigits	47.52(16)	99.99(2)	99.98(5)	99.99(2)	99.98(5)	99.98(5)	99.98(5)	94.53(15)	100(1)	99.98(5)	99.1(14)	99.76(13)	99.98(5)	99.98(5)	99.97(11)	99.97(11)
PageBlocks	79.45(16)	98.08(7)	95.25(12)	89.15(15)	90.1(14)	95.55(11)	98.73(4)	92.29(13)	96.88(9)	96.49(10)	97.61(8)	98.44(6)	98.78(3)	98.82(2)	98.62(5)	99.05(1)
Parkinson	99.05(8)	84.19(15)	91.21(12)	91.13(12)	95.05(11)	99.86(4)	83.34(14)	97.37(9)	97.27(9)	96.27(10)	98.8(3)	99.81(5)	99.67(7)	99.8(6)	99.92(1)	99.91(2)
pendigits	62.56(16)	99.99(2)	99.91(9)	99.79(14)	99.82(13)	99.91(9)	99.97(5)	99.62(15)	100(1)	99.95(7)	99.91(2)	99.97(5)	99.99(2)	99.95(7)	99.91(9)	99.99(2)
Pima	65.21(6)	83.81(10)	75.97(15)	83.25(13)	83.38(12)	82.79(14)	95.47(1)	83.47(11)	86.61(7)	84.51(9)	85.94(8)	88.07(6)	93.84(1)	92(3)	91.3(4)	92.02(2)
satellite	77.11(16)	96.22(5)	81.55(14)	84.27(13)	84.99(12)	86.32(11)	90.03(6)	78.66(15)	94.43(9)	91.78(10)	94.6(8)	95.88(7)	97.08(1)	96.59(3)	96.56(4)	96.98(2)
satellite-2	99.09(7)	97.1(4)	99.4(5)	99.4(5)	98.39(10)	99.74(1)	98.21(1)	98.14(12)	97.97(13)	98.73(8)	97.11(6)	97.57(15)	98.73(9)	99.1(6)	99.7(2)	99.72(1)
shuttle	78.8(16)	99.43(9)	98.69(1)	97.57(14)	97.59(13)	98(21)	99.99(3)	97.74(12)	98.99(7)	97.54(15)	98.75(8)	98.55(10)	100(1)	100(1)	99.99(3)	99.99(3)
skin	52.94(16)	99.92(8)	89.31(5)	95.77(12)	95.28(13)	98.77(11)	99.97(1)	93.9(14)	99.94(4)	99.96(2)	99.99(10)	99.93(6)	99.94(4)	99.93(6)	99.96(2)	99.96(2)
smtp	99.23(5)	91.98(8)	85.48(11)	78.88(13)	86.42(9)	98.76(6)	99.01(4)	84.13(12)	86.42(9)	100(1)	95.83(7)	99.99(2)	56.21(15)	99.7(4)	99.99(2)	99.99(2)
Spambase	99.25(16)	95.94(10)	93.21(3)	92.03(14)	94.01(12)	94.35(11)	97.88(3)	84.31(15)	96.61(7)	97.33(5)	96.35(8)	96.35(8)	97.28(6)	97.05(2)	97.58(2)	98.14(1)
speech	47.76(16)	63.41(14)	79.23(6)	84.01(2)	85.32(1)	75.16(10)	98.26(2)	74.49(15)	82.3(12)	81.29(4)	73.85(12)	77.29(9)	60.61(5)	69.39(13)	79.68(5)	77.38(7)
Stamps	76.83(16)	98.69(13)	98.71(2)	98.42(14)	99.08(10)	98.34(15)	99.88(2)	99.08(29)	99.21(8)	99.21(8)	99.67(7)	99.75(6)	99.86(4)	99.88(2)	99.85(5)	99.93(1)
thyroid	93.27(16)	99.55(15)	99.8(8)	99.82(7)	99.77(10)	99.8(8)	99.95(1)	99.65(13)	99.77(10)	99.72(12)	99.58(14)	99.86(6)	99.89(4)	99.89(4)	99.93(2)	99.92(3)
vertebral	40.75(16)	92.02(10)	81.76(14)	81.74(15)	82.32(13)	93.53(8)	99.14(3)	96.78(11)	92.61(9)	84.63(12)	97.88(7)	98.29(6)	99.25(2)	98.82(4)	98.82(4)	99.35(1)
vowels	54.25(16)	95.53(14)	99.26(7)	98.52(9)	98.81(7)	98.41(11)	98.59(12)	29.25(15)	92.72(2)	95.52(15)	99.1(12)	99.88(1)	97.49(12)	97.22(4)	99.22(5)	99.22(5)
Waveform	52.91(16)	90.76(12)	92.28(8)	92.43(7)	93.48(4)	90.76(12)	93.3(5)	90.38(14)	94.22(2)	96.39(1)	83.91(15)	91.53(10)	90.98(11)	93.12(6)	92.01(9)	93.7(3)
WBC	95.93(16)	98.93(13)	99.18(10)	99.73(4)	99.1(11)	99.38(8)	99.69(5)	98.91(14)	98.97(12)	98.85(15)	99.91(2)	99.63(6)	99.93(1)	99.4(7)	99.31(9)	99.82(3)
WDBC	97.85(16)	99.99(11)	100(1)	100(1)	100(1)	100(1)	99.52(15)	100(1)	100(1)	100(1)	100(1)	100(1)	99.99(11)	99.95(14)	99.96(13)	100(1)
Wilt	46.01(16)	98.85(9)	99.81(3)	68.15(14)	69.11(13)	92.10(10)	98.56(7)	82.67(11)	98.88(2)	82.84(12)	99.01(2)	98.71(6)	98.81(9)	98.81(9)	98.81(9)	98.81(9)
wine	74.38(16)	99.98(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
WPBC	49.93(16)	95.22(8)	80.56(14)	84.11(12)	84.08(13)	91.45(10)	98.83(2)	71.44(15)	91.87(9)	90.37(11)	96.87(7)	97.27(6)	96.88(3)	99.33(1)	98.37(5)	98.59(4)
yeast	49.05(15)	70.94(9)	46.84(16)	69.04(12)	68.99(13)	70.32(11)	74.72(5)	67.26(14)	71.09(8)	70.58(10)	71.47(7)	72.29(6)	75.82(4)	74.97(3)	74.75(4)	78.02(1)
FashionMNIST	80.23(15)	98.53(14)	96.64(9)	88.09(1)	88.81(7)	95.39(9)	99.99(3)	72.55(16)	94.54(10)	93.7(11)	91.81(4)	92.88(13)	96.71(6)	96.71(6)	96.71(6)	96.71(6)
CIFAR10	61.85(16)	80.15(11)	86.13(4)	86.6(2)	84.02(8)	81.32(10)	85.87(5)	66.03(15)	84.43(7)	81.41(9)	74.42(13)	74.91(12)	73.82(14)	86.27(3)	85.56(6)	86.81(1)
SVHN	57.63(15)	77.98(10)	82.41(2)	83.44(1)	82(3)	97.29(9)	81.23(8)	53.15(16)	81.3(7)	77.03(11)	75.33(12)	70.31(13)	68.77(14)	81.54(6)	81.63(5)	81.88(4)
Agnews	57.92(16)	92.42(10)	92.48(8)	92.06(11)	91.28(12)	92.44(9)	93.24(6)	69.18(15)	94.26(2)	95.9(1)	93.16(7)	90.01(13)	77.11(14)	93.32(5)	93.53(4)	93.7(3)
Amazon	42.51(16)	83.75(11)	87.46(6(5))	87.46(6)	87.46(6)	84.51(10)	86.26(2)	52.28(16)	86.63(2)	91.58(1)	83.51(15)	86.96(14)	81.51(15)	86.37(10)	86.37(10)	86.37(10)
Imdb	49.91(16)	81.71(2)	83.5(8)	83.37(10)	82.85(11)	83.96(6)	83.56(7)	61.03(15)	66.83(2)	89.39(1)	83.9(8)	80.53(13)	64.42(14)	84.38(4)	84.11(5)	85.43(3)
Yelp	67.99(15)	89.63(11)	93.15(3)	92.48(4)	90.78(10)	91.84(6)	91.74(7)	59.81(16)	94.49(2)	95.24(1)	89(13)	89.07(12)	70.94(14)	91.63(8)	91.61(9)	92.41(5)

Table 16: AUCPR of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 75\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	67.03(16)	83.92(5)	78.15(15)	81.53(8)	81.46(10)	79.88(14)	83.61(6)	80.36(12)	85.74(1)	85.55(2)	80.91(11)	85.27(3)	81.71(7)	81.52(9)	80.33(13)	81.19(4)
ALOI	3.9(16)	8.77(7)	4.49(12)	4.19(13)	4.17(15)	5.82(10)	16.99(5)	4.25(16)	9.13(6)	7.17(8)	6.56(19)	5.05(11)	24.07(1)	18.29(2)	18.05(4)	19.77(2)
anthyroid	40.97(16)	73.75(10)	42.93(15)	44.64(14)	45.04(13)	59.76(11)	88.78(5)	47.44(12)	79.51(7)	71.15(9)	80.98(7)	86.84(6)	90.61(1)	89.84(1)	90.59(2)	90.44(3)
breast	92.82(1)	93.13(1)	93.48(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)	93.87(1)
breastw	89.92(16)	98.75(10)	98.73(15)	98.73(15)	99.48(4)	98.96(12)	99.29(8)	99.48(4)	98.95(11)	99.53(8)	99.29(8)	99.29(8)	99.29(8)	99.29(8)	99.29(8)	99.29(8)
cardio	49.13(15)	91.99(9)	95.63(3)	91.61(1)	91.91(10)	91.06(13)	96.2(1)	79.33(14)	95.71(2)	91.16(2)	NA/NA	95.27(5)	95.36(4)	94.22(8)	94.87(7)	95.11(6)
Cardiography	90.43(15)	88.59(9)	88.93(10)	84.29(13)	85.47(12)	79.88(15)	92.45(8)	80.36(12)	89.52(7)	90.16(2)	87.42(11)	89.39(8)	93.58(4)	92.58(10)	94.23(5)	93.36(6)
concrete	85.92(13)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)	91.82(3)
cover	66.48(16)	80.5(12)	80.28(15)	88.17(1)	88.4(10)	85.61(3)	93.42(2)	81.78(14)	91.67(8)	92.93(3)	90.46(9)	92.57(7)	92.75(6)	92.82(5)	92.87(4)	95.31(1)
fault	0.89(16)	98.28(1)	96.36(8)	97.34(4)	96.83(6)	93.7(13)	93.06(14)	91.95(11)	97.82(2)	96.01(9)	97.35(3)	94.79(11)	97.03(5)	94.82(10)	92.76(12)	96.79(7)
glass	54.25(16)	69.03(9)	66.75(11)	62.4(14)	67.29(10)	63.61(13)	72.7(15)	55.44(15)	70.49(7)	72.67(6)	66.47(12)	68.94(9)	78.01(2)	76.73(7)	75.54(4)	78.59(1)
plant	77.61(13)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)	93.3(1)
HeartDisease	49.93(9)	93.93(8)	91.4(14)	91.43(3)	91.81(11)	90.67(15)	98.26(5)	91.56(12)	93.42(10)	94.38(8)	97.57(7)	97.87(6)	99.05(1)	98.64(1)	98.71(3)	98.71(3)
Hepatitis	37.56(16)	99.62(8)	99.96(11)	93.59(14)	96.74(12)	94.23(15)	99.59(9)	82.09(15)	99.34(10)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
hsp	73.48(14)	93.92(1)	93.92(1)	100(1)	100(1)	90.13(1)	96.68(9)	96.68(9)	100(1)	100(1)	96.68(10)	99.44(8)	99.44(8)	99.44(8)	99.44(8)	99.44(8)
Image	65.88(16)	88.89(13)	88.36(14)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)	91.9(11)
Internet Ads	51.82(14)	88.69(8)	91.96(3)	92.19(2)	91.55(4)	85.33(10)	89.37(7)	47.21(15)	91.45(8)	84.36(12)	84.81(11)	NA/NA	84.05(13)	87.59(9)	90.22(6)	92.58(1)
Ionosphere	91.39(15)	98.06(6)	97.97(9)	94.01(13)	93.95(14)	98.02(7)	99.4(1)	88.57(16)	97.81(11)	97.67(10)	97.07(12)	98.02(7)	99.07(3)	98.94(4)	99.12(2)	98.91(5)
landsat	24.48(16)	85.55(6)	42.47(14)	57.03(13)	57.88(12)	60.29(11)	87.66(4)	29.25(15)	82.06(9)	82.33(8)	78.22(10)	84.05(7)	89.42(2)	88.53(1)	87.56(5)	89.57(1)
letter	77.95(13)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)	94.41(3)
Lymphography	76.38(16)	86.1(7)	77.37(11)	73.85(14)	75.38(13)	77.37(11)	87.18(5)	69.45(15)	87.08(6)	87.58(4)	82.94(10)	83.75(9)	88.34(2)	87.71(3)	85.72(8)	88.99(1)
magic-gamma	14.43(16)	66.66(17)	62.40(10)	62.24(11)	61.19(12)	58.78(13)	66.55(6)	48.45(15)	56.34(14)	63.16(9)	68.39(4)	71.68(2)	69.46(3)	68.58(5)	66.92(5)	74.01(1)
mammography	16.41(16)	66.29(10)	65.45(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)	66.29(10)
mist	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
mnist	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
optdigits	3.23(16)	97.72(2)	99.6(5)	96.63(4)	99.67(3)	99.37(10)	99.38(8)	29.78(15)	100(1)	99.59(6)	97.47(14)	98.57(13)	99.38(1)	99.80(1)	99.45(7)	99.07(12)
PageBlocks	5.44(16)	87.97(8)	79.71(11)	67.37(14)	69.38(13)	74.23(12)	88.5(6)	55.98(15)	85.16(9)	84.1(10)	80.23(5)	89.84(3)	90.23(2)	89.47(4)	88.25(7)	91.76(1)
Pen-Kinship	99.66(8)	97.72(10)	97.05(12)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)	97.72(10)
pen-test	47.4(16)	99.7(3)	97.55(9)	93.21(5)	94.63(13)	96.28(12)	99.01(5)	94.11(14)	99.88(1)	98.72(8)	97.23(1)	99.72(1)	99.72(1)	98.97(1)	99.05(1)	99.44(1)
Pima	47.65(16)	69.97(14)	61.29(15)	70.81(12)	70.85(11)	73.46(10)	83.7(5)	70.78(13)	77.73(8)	75.04(9)	78.85(6)	78.85(6)	90.24(1)	84.13(2)	87.26(4)	88.08(3)
satellite	67.69(16)	92.78(6)	80.24(13)	82.61(12)	83.15(11)	80.19(14)	93.42(5)	75.71(15)	91.67(7)	87.99(10)	88.68(9)	95.85(4)	95.05(1)	88.31(3)	93.45(4)	94.9(2)
satimgm-2	50.1(16)	96.1(11)	93.81(7)	93.58(8)	91.96(11)	93.97(10)	95.49(2)	91.51(2)	91.51(2)	89.99(9)	89.32(15)	91.41(3)	89.71(4)	94.61(4)	94.85(3)	94.44(6)
satimgm-3	68.7(16)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)	97.1(10)
skin	22.52(16)	99.56(6)	49.15(15)	70.63(13)	67.45(14)	89.17(11)	99.86(1)	85.56(12)	99.33(8)	99.82(2)	99.15(10)	99.22(9)	99.68(4)	99.66(5)	99.57(7)	99.55(5)
smtp	20.87(15)	62.14(10)	66.71(5)	66.69(6)	66.68(8)	66.34(12)	50.31(10)	67.16(16)	66.69(6)	50.03(11)	100(1)	34.5(13)	83.33(2)	33.33(2)	33.33(2)	33.33(2)
SpamBase	45.24(15)	97.83(10)	88.75(13)	86.87(14)	90.32(12)	91.44(11)	96.89(3)	68.56(15)	95.47(6)	95.02(5)	93.36(8)	92.99(9)	96.63(4)	97.07(4)	97.07(4)	97.07(4)
statlog	9.82(17)	91.94(2)	88.27(7)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)	87.77(9)
Stamps	34.54(16)	80.17(14)	85.33(11)	77.27(15)	82.96(12)	82.52(13)	99.21(2)	88.76(10)	94.68(8)	90.9(9)	96.32(7)	97.46(6)	96.66(4)	96.76(3)	98.49(5)	99.38(1)
thyroid	15.55(16)	90.73(12)	92.99(9)	94.4(10)	94.4(10)	93.56(8)	98.29(1)	90.22(11)	92.22(11)	92.01(12)	91.82(13)	96.07(4)	96.96(5)	98.28(5)	95.15(6)	97.16(2)
vertebral	55(16)	69.15(9)	60.63(13)	46.16(15)	46.16(15)	46.16(15)	46.16(15)	52.55(12)	74.78(8)	55.19(11)	67.83(6)	90.99(3)	90.99(3)	90.99(3)	90.99(3)	90.99(3)
Waveform	28.43(16)	97.72(10)	92.72(15)	91.77(12)	91.77(12)	90.29(14)	99.86(1)	89.29(14)	99.86(1)	99.86(1)	99.86(1)	99.86(1)	99.86(1)	99.86(1)	99.86(1)	99.86(1)
Waveform	47.47(16)	47.7(6)	19.72(15)	19.83(14)	24.14(11)	23.36(12)	51.91(3)	22.56(13)	66.56(1)	53.82(2)	40.49(8)	43.72(9)	50.61(5)	46.28(7)	51.39(4)	51.39(4)
WBC	48.77(16)	89.62(12)	89.93(11)	91.77(9)	84.18(15)	95.14(7)	95.85(1)	88.35(13)	94.55(10)	84.55(14)	98.1(2)	95.9(4)	98.46(1)	95.91(6)	95.98(1)	97.04(3)
WDBC	61.66(16)	99.63(11)	100(1)	100(1)	100(1)	100(1)	96.98(15)	100(1)	100(1)	100(1)	100(1)	100(1)	99.54(12)	98.96(14)	99.07(13)	100(1)
wine	84.42(9)	84.42(9)	81.1(13)	81.94(14)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)	81.94(13)
wine	24.84(16)	99.85(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
WPBC	23.91(16)	86.09(9)	64.58(14)	69.31(13)	70.04(12)	84.92(10)	97.86(5)	48.03(15)	87.85(8)	81.52(11)	96.25(7)	97.43(6)	98.11(2)	98.76(1)	98.04(3)	98.04(3)
yeast	33.54(16)	51.46(1)	36.36(15)	49.24(12)	48.99(13)	54.83(7)	60.76(5)	46.92(14)	52.32(9)	51.54(10)	55.07(6)	54.38(8)	63.83(2)	61.78(3)	61.63(4)	65.46(4)
YeastMNIST	22.08(16)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)
CIFAR10	9.42(16)	31.81(1)	41.15(1)	38.42(4)	36.93(8)	33.51(9)	38.04(7)	9.59(15)	38.16(6)	32.83(10)	30.76(12)	24.34(13)	19.97(14)	39.12(5)	39.26(5)	39.37(2)
SVHN	7.89(15)	70.02(11)	34(3)	34.17(2)	34.87(1)	27.04(12)	31.61(8)	5.39(16)	33.62(4)	30.71(9)	29.51(10)	20.43(13)	24.14(14)	33.54(5)	32.96(6)	32.53(7)
Agnews	2.99(16)	20.91(4)	64.27(5)	55.82(12)	58.75(11)	61.56(9)	59.97(10)	8.77(15)	71.95(2)	74.02(10)	70.14(3)	55.13(13)	25.48(14)	62.2(7)	61.82(8)	64.27(5)
Agnews	35.79(15)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)	92.3(10)
Imdb	5(16)	34.28(4)	29.92(8)	27.08(11)	26.51(12)	28.58(10)	28.68(9)	7.75(15)	38.36(2)	44.09(1)	35.61(3)	21.04(13)	13.27(14)	30.36(7)	31.13(6)	31.46(5)
Yelp	9.02(15)	56.25(3)	53.63(4)	49.62(7)	50.44(5)	46.06(11)	46.38(9)	65.21(16)	64.11(11)	63.68(2)	45.32(12)	41.14(13)	13.27(14)	47.41(8)	46.2(10)	50.22(5)

Table 17: AUCROC of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 100\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	72.33(16)	86.83(2)	77.79(15)	82.1(13)	82.48(12)	85.15(10)	85.94(6)	81.41(14)	86.3(4)	86.17(5)	84.31(11)	86.74(3)	85.93(7)	85.6(8)	85.49(9)	86.85(1)
ALOI	57.04(12)	68.95(6)	52.22(16)	53.84(14)	53.32(15)	65.82(7)	82.38(2)	54.35(13)	62.88(9)	64.54(8)	60.41(0)	57.47(11)	82.85(1)	79.37(3)	76.55(5)	78.54(4)
anthyroid	62.58(16)	88.17(10)	84.54(13)	82.7(15)	82.81(14)	89.96(12)	99.57(2)	91.67(11)	98.72(8)	98.62(9)	98.98(7)	99.13(6)	99.6(1)	99.56(4)	99.57(2)	99.55(5)
Arrhythmia	68.56(16)	96.83(7)	93(10)	88(13)	88.15(12)	90.1(11)	98.17(5)	85.15(15)	97.15(14)	94.04(9)	97.31(6)	96.67(8)	99.43(1)	98.38(3)	98.58(4)	99.01(2)
breastw	97.28(16)	99.82(8)	99.01(14)	99.72(11)	99.77(10)	99.1(13)	99.9(5)	98.87(15)	99.66(12)	99.82(8)	99.95(1)	99.86(7)	99.92(3)	99.92(3)	99.87(6)	99.93(2)
cardio	87.96(15)	99.2(10)	99.42(9)	98.96(11)	98.87(12)	98.21(13)	99.81(2)	97.64(14)	99.71(4)	99.44(8)	NANA	99.82(1)	99.73(3)	99.64(5)	99.49(7)	99.61(6)
Cardiotocography	64.62(16)	96.95(9)	96.09(11)	94.82(13)	95.38(12)	94.44(14)	98.14(6)	93.75(15)	96.87(10)	97.21(8)	98.01(7)	98.2(5)	98.33(3)	98.54(2)	98.28(4)	98.57(1)
comm.and.crimne	79.86(16)	89.41(2)	90.62(10)	91.63(4)	91.33(8)	86.46(15)	91.47(6)	88.59(14)	91.83(3)	91.36(2)	98.68(13)	90.26(11)	91.54(5)	91.47(6)	91.08(9)	92.26(1)
concrete	66.94(16)	89.77(12)	87.51(14)	90.4(10)	89.98(11)	89.46(13)	95.75(7)	85.61(5)	93.24(8)	93.12(9)	95.76(6)	95.91(5)	96.26(4)	96.8(2)	96.58(3)	97.11(1)
cover	46.16(16)	99.97(7)	99.95(11)	99.96(9)	99.96(9)	99.91(13)	99.89(15)	99.94(14)	99.98(3)	99.97(7)	99.99(1)	99.99(1)	99.98(3)	99.98(3)	99.93(12)	99.98(3)
fault	83.57(8)	78.55(12)	77.29(13)	79.21(11)	76.51(14)	85.29(6)	70.44(15)	82.56(10)	84.49(7)	83.31(9)	87.41(5)	89.22(3)	89.64(2)	88.65(4)	89.8(1)	
glass	68.91(6)	98.51(10)	89.45(14)	89.78(15)	90.44(13)	100(1)	99.98(7)	93.71(1)	99.47(9)	92.67(12)	99.73(8)	100(1)	100(1)	100(1)	100(1)	100(1)
HeartDisease	76.26(16)	95.98(8)	92.37(15)	92.99(13)	93.21(12)	95.44(11)	99.55(2)	92.62(14)	95.63(10)	95.77(9)	99.27(5)	99.06(6)	99.86(1)	99.52(3)	98.93(7)	99.35(4)
Hepatitis	75.9(16)	100(1)	99.72(12)	99.09(14)	99.45(13)	99.93(10)	99.85(11)	97.41(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
http	99.81(5)	100(1)	99.98(13)	100(1)	100(1)	99.84(14)	100(1)	98.33(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
imgsg	82.23(15)	98.26(8)	87.52(13)	88.29(12)	87.49(14)	95.21(11)	99.48(7)	81.03(16)	96.44(10)	96.53(9)	99.71(6)	99.75(5)	99.84(4)	99.87(3)	99.88(2)	99.9(1)
InternetAds	70.14(15)	96.5(7)	96.49(8)	97.11(3)	97.16(2)	94.11(3)	95.75(9)	85.96(14)	95.28(11)	97.63(1)	97(5)	NANA	97.04(4)	95.14(12)	95.37(10)	96.91(6)
Ionosphere	93.9(16)	99.25(8)	98.29(11)	94.76(14)	94.88(13)	98.12(12)	100(1)	93.96(15)	98.92(10)	99.07(9)	99.98(5)	99.56(7)	100(1)	99.98(5)	100(1)	99.99(4)
landsat	57.87(16)	95.62(8)	60.98(15)	80.41(12)	80.31(13)	86.86(11)	95.89(7)	71.88(14)	94.77(9)	94.52(10)	96.35(5)	96.27(6)	96.62(4)	96.78(3)	96.79(2)	96.8(1)
letter	70.88(16)	87.01(12)	90.09(9)	87.71(1)	86.17(13)	84.35(14)	94.93(4)	75.85(15)	90.82(8)	89.59(10)	94.85(3)	91.41(7)	96.64(2)	90.2(3)	94.02(6)	97.3(1)
Lymphography	98.4(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	99.59(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
magic.gamma	56.75(16)	91.45(8)	83.17(14)	83.36(12)	83.31(13)	86.11(11)	91.68(6)	75.84(15)	90.61(10)	91.17(9)	91.67(7)	91.95(5)	92.49(3)	92.68(2)	92.37(4)	93.15(1)
mammography	76.95(16)	95.75(3)	92.88(14)	93.14(12)	93.24(11)	95.84(2)	95.64(4)	92.96(13)	87.43(15)	95.17(5)	93.85(8)	95.06(6)	94.72(7)	93.49(10)	93.63(9)	95.95(1)
mnist	72.74(16)	99.68(8)	99.29(11)	99.28(12)	99.25(12)	98.08(14)	97.73(6)	95.35(16)	99.69(7)	99.77(5)	99.89(10)	99.89(9)	99.83(1)	99.82(2)	99.83(3)	99.85(1)
musk	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
optdigits	51.61(16)	100(1)	99.98(10)	99.98(10)	99.99(7)	99.97(13)	99.99(7)	94.31(13)	99.97(13)	100(1)	100(1)	100(1)	99.98(10)	100(1)	99.99(7)	100(1)
PageBlocks	79.86(16)	98.51(8)	95.66(12)	89.24(15)	90.83(14)	96(11)	99.01(5)	93.13(15)	96.88(10)	97.29(9)	98.53(7)	98.72(6)	99.02(4)	99.25(2)	99.13(3)	99.31(1)
Parkinson	76.94(16)	99.82(8)	89.13(14)	92.99(12)	92.53(13)	98.18(9)	100(1)	88.18(15)	97.29(10)	96.31(11)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
pendigits	57.49(16)	99.99(2)	99.88(10)	99.79(12)	99.77(14)	99.84(11)	99.97(7)	99.61(15)	100(1)	99.94(8)	99.99(2)	99.98(4)	99.98(4)	99.93(9)	99.78(13)	99.98(4)
Pima	66.35(16)	85.52(10)	77.22(15)	83.74(11)	83.44(13)	83.13(14)	92.44(5)	83.66(12)	87.12(8)	85.59(9)	88.26(7)	90.29(6)	96.12(1)	93.96(3)	93.32(4)	94.17(2)
satellite	80.79(15)	96.74(7)	81.89(14)	83.58(13)	83.96(12)	88.24(11)	96.72(8)	78.28(16)	95.53(9)	94.38(10)	96.87(6)	97.32(5)	97.41(3)	97.41(3)	97.51(1)	97.49(2)
satellite-2	99.57(11)	99.68(2)	99.53(5)	99.5(4)	99.47(5)	98.19(15)	99.22(6)	98.51(2)	99.17(8)	97.94(16)	96.88(8)	98.82(10)	98.33(14)	98.42(13)	98.87(9)	99.76(1)
shuttle	96.34(16)	99.52(7)	97.57(15)	98.58(14)	97.58(14)	97.82(15)	100(1)	99.13(8)	99.02(10)	98.21(2)	99.71(6)	99.13(8)	100(1)	100(1)	100(1)	100(1)
skin	51.02(16)	99.84(10)	88.75(15)	96.19(12)	95.34(13)	98.71(11)	99.99(1)	93.93(14)	99.95(9)	99.96(7)	99.98(3)	99.97(6)	99.99(1)	99.98(3)	99.96(7)	99.98(3)
smtp	55.31(16)	99.23(5)	91.98(8)	85.48(10)	78.88(13)	85.3(11)	98.76(6)	59.01(14)	84.13(12)	86.42(9)	100(1)	95.83(7)	99.99(2)	96.21(15)	99.7(4)	99.99(2)
Spambase	97.29(9)	94.74(12)	92.29(14)	94.11(13)	96.37(11)	98.18(4)	98.14(6)	86.42(15)	97.12(10)	97.88(7)	97.68(8)	98.01(5)	98.51(11)	98.43(3)	98.51(1)	
speech	65.84(14)	81.01(5)	99.03(9)	99.03(9)	98.92(10)	99.29(8)	99.29(8)	97.86(14)	99.98(1)	96.29(15)	97.47(11)	96.68(4)	95.62(15)	74.82(12)	80.96(5)	78.78(8)
Stamps	78.11(16)	99.05(12)	98.97(13)	98.54(15)	99.09(11)	99.31(9)	99.98(3)	98.97(13)	99.65(8)	99.23(10)	99.96(4)	99.88(7)	99.99(3)	100(1)	99.89(6)	100(1)
thyroid	93.19(16)	99.61(15)	99.85(8)	99.81(9)	99.79(10)	99.76(13)	99.96(3)	99.69(14)	99.78(11)	99.78(11)	99.87(7)	99.93(5)	99.97(1)	99.93(5)	99.96(3)	99.97(1)
vertebral	43.13(16)	93.38(8)	82.43(13)	81.08(15)	82.18(14)	93.04(9)	99.64(1)	86.48(12)	91.5(10)	86.88(11)	99.25(5)	99.01(7)	99.43(2)	99.02(6)	99.34(4)	99.43(2)
vowels	98.07(16)	98.07(13)	99.45(7)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)	99.03(9)
Waveform	53(16)	93.02(13)	93.21(11)	93.41(0)	93.97(5)	93.15(12)	93.75(7)	91.21(14)	95.36(2)	96.84(1)	81.72(15)	93.85(6)	93.53(9)	94.81(4)	93.62(8)	94.86(5)
WBC	96.14(16)	98.92(12)	99.17(11)	99.56(10)	98.88(14)	99.71(9)	100(1)	98.55(15)	100(1)	98.92(12)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
WDBC	98(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
Wilt	58.63(16)	98.96(8)	58.15(15)	68.16(14)	68.68(13)	64.69(11)	97.14(2)	83.61(12)	96.86(9)	98.68(10)	99.14(2)	98.86(7)	98.89(9)	99.07(4)	99(5)	
wine	80.09(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
WPBC	51.56(16)	96.79(8)	81.07(14)	85.62(12)	84.19(13)	94.44(10)	99.93(5)	72.04(15)	95.34(9)	92.14(11)	99.47(7)	99.88(6)	100(1)	100(1)	99.96(4)	99.98(3)
yeast	49.03(15)	72.41(8)	46.92(16)	67.54(13)	67.95(12)	68.17(11)	75.95(5)	71.7(9)	71.15(10)	73.88(7)	75.64(6)	78.36(2)	77.69(3)	77.11(4)	79.11(4)	79.41(1)
FashionMNIST	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)	96.64(15)
CIFAR10	62.48(16)	82.71(11)	86.92(5)	87.4(2)	85.33(8)	82.92(10)	87.37(3)	69.33(15)	86.12(7)	83.03(9)	78.88(12)	77.5(13)	75.75(14)	87.32(4)	86.61(6)	87.63(1)
SVHN	58.26(15)	80.02(9)	83.13(3)	84.16(1)	83.4(2)	78.62(12)	81.98(8)	54.97(16)	82.57(7)	78.98(11)	79.21(0)	72.36(13)	70.44(14)	82.71(5)	82.72(4)	82.59(6)
Agnews	57.95(16)	94(5)	93.38(8)	92.37(11)	92.31(12)	92.39(10)	93.1(9)	68.87(15)	94.56(3)	96.49(1)	95.24(2)	91.87(13)	80.81(4)	93.87(7)	93.98(6)	94.04(4)
Amazon	50.66(15)	83.14(12)	82.83(10)	83.79(11)	83.79(11)	84.5(8)	84.66(7)	66.93(14)	82.16(16)	91.5(2)	87.51(12)	84.66(4)	85.12(13)	97.99(1)	87.62(6)	88.81(1)
Imdb	50.05(16)	81.07(12)	82.83(10)	82.79(11)	84.5(8)	84.66(7)	66.93(14)	66.93(14)	82.16(16)	91.5(2)	87.51(12)	84.66(4)	85.12(13)	97.99(1)	87.62(6)	88.81(1)
Yelp	68.9(15)	91.06(11)	93.29(3)	92.18(6)	92.05(9)	91.19(10)	92.18(6)	60.09(16)	95.29(2)	95.65(1)	90.61(2)	89.17(13)	76.25(14)	92.65(5)	92.12(8)	92.77(4)

Table 18: AUCPR of 16 label-informed algorithms on 55 real-world datasets, with labeled anomaly ration $\gamma_l = 100\%$. We show the performance rank in parenthesis (lower the better).

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTrans former	RF	LGB	XGB	CatB
abalone	71.75(16)	85.83(2)	78.49(15)	81.43(13)	81.88(12)	84.92(6)	85(5)	80.57(14)	84.86(7)	85.73(4)	82.77(11)	85.82(3)	84.25(8)	84.23(9)	83.93(10)	86.15(1)
ALOI	3.92(16)	9.31(7)	4.57(12)	4.31(3)	4.19(15)	6.02(10)	18.12(5)	42.14(1)	41.46(6)	9.06(8)	7.68(9)	5.01(11)	31.35(1)	26.05(2)	23.15(4)	25.44(3)
anthyroid	45.77(13)	78.04(10)	45.07(15)	45.35(14)	44.64(16)	63.95(12)	93.23(2)	60.64(11)	80.97(8)	79.61(9)	85.09(7)	86.5(6)	93.25(1)	92.97(3)	92.44(4)	92.36(5)
Arrhythmia	69.77(16)	85.83(2)	78.49(15)	81.43(13)	81.88(12)	84.92(6)	85(5)	80.57(14)	84.86(7)	85.73(4)	82.77(11)	85.82(3)	84.25(8)	84.23(9)	83.93(10)	86.15(1)
breastw	94.78(16)	99.66(9)	97.83(4)	99.46(1)	99.56(10)	98.99(13)	99.81(5)	96.36(11)	99.22(12)	99.67(8)	99.9(1)	99.71(7)	99.93(4)	99.87(5)	99.74(6)	99.82(5)
cardio	53.07(15)	95.14(10)	98.27(9)	92.91(13)	93.03(12)	93.16(1)	98.46(2)	81.01(14)	87.85(3)	96.48(8)	NANA	98.65(1)	97.79(4)	97.41(5)	96.96(7)	97.46(1)
Cardiography on normal and crine	44.34(16)	90.49(10)	96.79(11)	82.95(14)	85.97(12)	84.61(13)	94.35(8)	79.15(8)	89.57(10)	91.69(8)	93.19(7)	94.66(4)	95.12(3)	95.21(2)	93.56(6)	95.38(1)
concrete	66.9(16)	81.16(12)	85.77(14)	88.66(10)	88.57(11)	87.58(13)	95.48(6)	82.49(15)	82.85(8)	91.81(9)	95.31(7)	95.61(5)	95.85(4)	96.33(2)	96.21(3)	96.6(1)
cover	0.92(16)	19.18(6)	96.54(3)	97.44(9)	97.25(12)	94.12(14)	97.63(8)	92.12(15)	98.27(5)	97.36(11)	98.75(2)	98.95(1)	98.18(6)	98.37(4)	97.41(10)	98.45(3)
fault	55.63(16)	73.77(10)	69.59(11)	64.79(13)	66.92(12)	64.77(14)	77.52(6)	57.04(15)	74.86(9)	75.14(8)	76.6(7)	76.95(5)	82.69(3)	83.41(2)	81.64(4)	83.41(2)
glass	68.14(16)	88.14(12)	26.98(15)	88.92(10)	88.92(10)	99.36(9)	99.36(9)	83.74(11)	83.74(11)	83.74(11)	83.74(11)	83.74(11)	83.74(11)	83.74(11)	83.74(11)	83.74(11)
HeartDisease	70.61(16)	94.89(10)	91.41(1)	91.99(13)	92.08(12)	93.64(11)	99.45(3)	91.19(15)	94.90(9)	95.19(8)	99.39(5)	99.03(7)	99.81(1)	99.57(2)	99.09(6)	99.74(4)
Hepatitis	43.83(16)	100(1)	98.76(12)	95.28(14)	97.23(13)	99.60(10)	99.22(11)	82.15(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
img	73.48(16)	91.78(10)	98.65(12)	95.28(14)	97.23(13)	99.60(10)	99.22(11)	82.15(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
Internet Ads	73.48(16)	91.78(10)	98.65(12)	95.28(14)	97.23(13)	99.60(10)	99.22(11)	82.15(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
Ionosphere	54.44(14)	93.88(2)	92.42(9)	92.53(7)	93.31(5)	88.96(13)	90.96(11)	50.63(15)	93.07(6)	93.74(3)	94.09(1)	NANA	92.48(8)	90.97(12)	91.44(10)	93.62(4)
landsat	91.77(16)	99.02(8)	97.85(12)	94.06(14)	94.22(13)	98.37(11)	99.9(1)	92.82(15)	98.81(9)	98.81(9)	99.96(6)	99.45(7)	99.99(1)	99.87(5)	99.44(1)	99.98(4)
landsat	45.77(16)	86.53(8)	43.41(4)	56.64(13)	58.12(12)	62.71(1)	88.7(7)	29.25(15)	85.3(9)	83.89(10)	95.85(5)	90.25(4)	90.78(3)	90.82(1)	90.82(1)	90.82(1)
libras	21.58(16)	54.56(12)	32.36(13)	54.56(12)	54.56(12)	54.56(12)	54.56(12)	19.58(15)	74.7(2)	74.7(2)	74.7(2)	74.7(2)	74.7(2)	74.7(2)	74.7(2)	74.7(2)
Lymphography	78.35(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	85.01(15)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
magicgamma	52.2(16)	88.07(9)	77.59(12)	74.86(14)	75.47(13)	80.79(11)	88.68(7)	68.16(15)	87.51(10)	88.18(8)	88.77(6)	89.02(5)	89.37(4)	90.06(2)	89.68(3)	90.68(1)
mammography	45.77(16)	86.99(8)	61.99(12)	61.73(13)	61.56(14)	62.71(11)	70.55(6)	49.19(15)	63.45(10)	65.35(9)	70.2(7)	71.21(5)	73.71(2)	74.94(2)	72.21(4)	76.8(1)
meas	100(1)	100(1)	97.1(11)	90.6(13)	92.54(12)	93.54(10)	79.8(8)	86.6(15)	98.7(8)	98.1(8)	98.1(8)	98.1(8)	98.1(8)	98.1(8)	98.1(8)	98.1(8)
musk	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
optdigits	3.71(16)	99.86(5)	99.57(11)	99.61(10)	99.7(8)	99.44(13)	99.7(8)	21.13(15)	100(1)	99.43(14)	100(1)	99.92(3)	99.51(2)	99.89(4)	99.75(7)	99.86(1)
PageBlocks	52.66(16)	90.21(8)	80.32(11)	67.63(14)	70.75(13)	75.78(12)	92.15(7)	60.98(15)	84.8(10)	86.22(9)	91.69(7)	92.34(4)	93.01(2)	92.85(5)	91.75(6)	93.54(1)
pen-digits	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)	99.94(16)
pendigits	4.99(16)	99.52(5)	97.37(11)	93.63(13)	93.54(14)	92.54(15)	99.34(6)	93.77(12)	99.80(1)	98.54(9)	97.67(2)	98.75(9)	99.56(4)	99.62(2)	99.16(8)	99.58(3)
Pima	48.61(16)	73.74(10)	63.61(5)	71.68(13)	72.62(11)	72.46(12)	88.71(5)	70.58(14)	78.43(8)	76.53(9)	84.89(6)	84.51(7)	93.79(1)	89.96(3)	89.51(4)	91.85(2)
satellite	75.97(15)	94.28(1)	81.46(4)	82.14(13)	82.58(12)	85.61(1)	94.59(7)	75.78(16)	92.89(9)	91.15(10)	94.97(6)	95.49(5)	95.95(4)	95.76(3)	95.86(1)	95.82(2)
shuttle	96.56(16)	99.83(10)	97.15(13)	97.15(13)	97.15(13)	97.15(13)	97.15(13)	96.56(15)	99.83(9)	99.83(9)	99.83(9)	99.83(9)	99.83(9)	99.83(9)	99.83(9)	99.83(9)
shuttle	82.87(16)	98.98(7)	97.19(10)	96.49(13)	96.43(14)	96.32(15)	99.93(3)	97.19(10)	97.76(9)	96.66(12)	99.28(6)	98.03(8)	100(1)	99.99(3)	99.99(3)	100(1)
skin	22.52(16)	99.18(10)	48.47(15)	73.57(13)	67.63(14)	99.96(11)	99.95(2)	84.26(12)	99.96(9)	99.96(9)	99.95(9)	99.82(6)	99.97(1)	99.94(3)	99.79(7)	99.92(4)
smtp	20.87(15)	67.14(14)	66.71(15)	66.69(6)	66.68(9)	66.69(6)	50.31(1)	70.16(16)	66.69(6)	50.03(12)	100(1)	34.5(13)	33.33(2)	33.37(14)	50.79(10)	33.33(2)
Spoken	95.67(16)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	95.67(16)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	97.15(9)	97.15(9)
speech	1.63(16)	61.91(5)	39.96(2)	20.85(8)	19.82(11)	13.53(13)	16.39(12)	21.23(7)	44.86(1)	24.7(4)	35.8(3)	23.71(5)	7.25(14)	20.1(10)	20.52(9)	22.34(6)
Stamps	36.99(16)	97.44(12)	86.72(13)	78.69(15)	84.96(14)	89.77(10)	99.18(5)	87.05(11)	97.34(8)	90.39(9)	99.64(4)	98.73(7)	99.83(1)	100(1)	99.86(10)	100(1)
thyroid	55.72(16)	81.99(14)	84.28(13)	93.83(9)	93.03(11)	92.31(12)	98.17(3)	89.67(15)	92.07(13)	93.22(10)	95.54(6)	94.99(7)	98.52(1)	97.41(5)	98.16(4)	98.68(1)
Vowel	72.7(16)	97.14(10)	45.36(15)	97.14(10)	97.14(10)	97.14(10)	97.14(10)	72.7(16)	97.14(10)	97.14(10)	97.14(10)	97.14(10)	97.14(10)	97.14(10)	97.14(10)	97.14(10)
vowels	28.96(16)	84.53(13)	92.75(7)	91.67(9)	91.72(8)	91.67(9)	85.22(12)	69.26(15)	99.53(1)	71.47(14)	99.44(2)	97.99(3)	87.46(11)	94.27(5)	93.02(6)	96.14(4)
Waveform	4.55(16)	52.55(8)	22.34(13)	21.5(4)	24.61(12)	32.51(11)	54.97(6)	20.41(15)	69.28(1)	61.74(2)	50.48(10)	56.22(5)	51.9(9)	56.86(4)	53.89(7)	58.13(3)
WBC	50.43(16)	85.45(13)	88.65(11)	92.76(10)	82.24(14)	95.94(1)	100(1)	74.65(15)	100(1)	85.59(12)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
WBC	50.43(16)	85.45(13)	88.65(11)	92.76(10)	82.24(14)	95.94(1)	100(1)	74.65(15)	100(1)	85.59(12)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
Wilt	4.93(16)	85.13(9)	6.71(5)	8.18(4)	3.56(13)	3.97(11)	93.01(4)	40.92(10)	85.15(8)	34.94(12)	94.38(2)	94.53(1)	90.59(7)	91.83(5)	90.86(6)	93.39(3)
wine	34.98(16)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)	100(1)
wine	24.72(16)	90.22(9)	67.49(14)	72.16(12)	68.99(13)	86.38(10)	99.84(5)	48.38(15)	92.87(8)	85.23(11)	99.13(7)	99.78(6)	100(1)	100(1)	99.89(4)	99.95(3)
WV	51.48(16)	96.11(11)	47.63(12)	92.1(1)	92.1(1)	92.1(1)	92.1(1)	51.48(16)	96.11(11)	47.63(12)	92.1(1)	92.1(1)	92.1(1)	92.1(1)	92.1(1)	92.1(1)
FashionMNIST	23.88(16)	86.7(8)	44.06(6)	84.07(9)	83.08(10)	76.64(14)	85.42(7)	29.05(15)	81.81(11)	85.08(8)	86.91(4)	89.17(2)	76.87(13)	87.11(1)	86.92(1)	86.29(5)
CIFAR10	9.58(16)	38.61(9)	49.09(1)	39.52(7)	39.39(8)	31.00(12)	40.66(1)	10.77(15)	42.52(1)	38(11)	38.22(10)	28.13(13)	23.98(14)	42.03(4)	41.35(5)	42.85(3)
SVHN	8.06(15)	31.84(11)	35.7(4)	36.07(7)	37.1(1)	24.98(12)	32.39(10)	6.16(16)	35.97(3)	34.03(9)	35.5(12)	63.87(13)	17.52(14)	34.16(7)	34.14(8)	34.56(6)
Agnews	67.56(16)	75.56(3)	65.79(6)	56.15(13)	57.31(12)	66.46(10)	61.84(11)	8.63(15)	72.66(4)	78.37(1)	76.52(5)	64.09(10)	64.65(8)	64.83(7)	66.27(5)	67.27(4)
Agnews	67.56(16)	75.56(3)	65.79(6)	56.15(13)	57.31(12)	66.46(10)	61.84(11)	8.63(15)	72.66(4)	78.37(1)	76.52(5)	64.09(10)	64.65(8)	64.83(7)	66.27(5)	67.27(4)
Imdb	5.05(16)	35.68(4)	29.74(9)	26.93(12)	27.52(11)	31.85(7)	28.82(11)	3.44(15)	48.29(1)	46.38(2)	40.71(3)	23.94(13)	11.26(14)	32.99(6)	48.83(8)	32.89(5)
Yelp	9.25(15)	59.64(3)	54.23(8)	48.89(9)	49.49(8)	42.3(13)	48.26(11)	6.53(16)	66(1)	65.54(2)	54.82(4)	42.96(12)	16.24(14)	30.38(7)	30.52(10)	51.38(1)