
ADBench: Anomaly Detection Benchmark

Songqiao Han^{1,*}, Xiyang Hu^{2,*}, Hailiang Huang^{1,*†}, Minqi 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 57 benchmark datasets, named ADBench. Our extensive experiments (98,436 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 efficiently 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 [94], rare disease detection [196], social media analysis [186] [193], and intrusion detection [88]. AD algorithms aim to identify data instances that deviate significantly from the majority of data objects [59] [139] [146] [160], and numerous methods have been developed in the last few decades [3] [85] [102] [103] [129] [156] [172] [198]. Among them, the majority are designed for tabular data (i.e., no time dependency and graph structure). Thus, we focus on the *tabular* AD algorithms and datasets in this work.

Although there are already some benchmark and evaluation works for tabular AD [25] [38] [42] [53] [166], 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 concerning anomaly types (e.g., local vs. global); (*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 CV and NLP 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 the 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 four types of anomalies; and (*iii*) algorithm robustness and stability under three settings of data corruptions. Fig. 1 provides an overview of ADBench.

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

* All authors contribute equally. Names are listed in alphabetical ordering by the last name.

† Corresponding authors. Direct technical questions to Minqi Jiang and Yue Zhao.

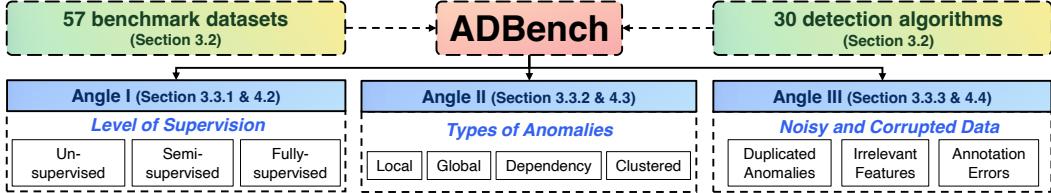


Figure 1: The design of the proposed ADBench is driven by research and application needs.

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 the 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 57 benchmark datasets (of which 47 are existing ones and we create 10).
2. **Research and application-driven benchmark angles.** By analyzing the needs of research and real-world applications, we focus on three critical comparison angles: availability of supervision, anomaly types, and algorithm robustness under noise and data corruption.
3. **Insights and future directions for researchers and practitioners.** With 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 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 [3]*, e.g., anomalies located in low-density regions, and their performance depends on the agreement between the input data and the algorithm assumption(s). Many unsupervised methods have been proposed in the last few decades [3] [15] [129] [150] [198], 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 [3], and surveys [129] [150] 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 [3] [170] such as Random Forest [21] and neural networks [89]. One known risk of supervised methods is that ground truth labels are not necessarily sufficient to capture all types of anomalies during annotation. These methods are therefore limited to detecting unknown types of anomalies [3]. Recent machine learning books [4] [54] and scikit-learn [133] may serve as good sources of supervised ML methods.

Semi-supervised Methods with Efficient Use of Labels. *Semi-supervised AD algorithms can capitalize the supervision from partial labels, while keeping the ability to detect unseen types of anomalies.* To this end, some recent studies investigate using partially labeled data for improving detection performance and leveraging 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 [7] [8] [188]. In ADBench, semi-supervision mostly refers to *incomplete label learning* in weak-supervision (see [205]). More discussions on semi-supervised AD are deferred to Appx. §A.3

2.2 Existing Datasets and Benchmarks for Tabular AD

AD Datasets in Literature. Existing benchmarks mainly evaluate a part of the datasets derived from the ODDS Library [145], DAMI Repository [25], ADRepository [129], and Anomaly Detection

Table 1: Comparison among ADBench and existing benchmarks, where ADBench comprehensively includes the most datasets and algorithms, uses both benchmark and synthetic datasets, covers both shallow and deep learning (DL) algorithms, and considers 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. [150]	3	9	✓	✓	✓	✓	✗	✓	✗
Goldstein et al. [53]	10	19	✓	✗	✓	✗	✗	✓	✗
Domingues et al. [38]	15	14	✓	✗	✓	✗	✗	✗	✓
Soenen et al. [164]	16	6	✓	✗	✓	✗	✗	✗	✗
Steinbuss et al. [166]	19	4	✗	✓	✓	✗	✗	✓	✗
Emmott et al. [42]	19	8	✓	✓	✓	✗	✗	✓	✓
Campos et al. [25]	23	12	✓	✗	✓	✗	✗	✗	✗
ADBench (ours)	57	30	✓	✓	✓	✓	✓	✓	✓

Meta-Analysis Benchmarks [42]. In ADBench, we include almost all publicly available datasets, and add larger datasets adapted from CV and NLP domains, for a more holistic view. See details in §3.2.

Existing Benchmarks. There are some notable works that take effort to benchmark AD methods on tabular data, e.g., [25] [38] [42] [150] [166] (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 that a comprehensive benchmark should also consider them. Second, most existing works only evaluate public benchmark 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 between the proposed ADBench and existing AD benchmarks in Table 1.

Also, “anomaly detection” is an overloaded term; there are AD benchmarks for time-series [85] [87] [132], graph [101], CV [6] [27] [202] and NLP [143], but they are different from tabular AD in nature.

2.3 Connections with Related Fields and Other Opportunities

While ADBench focuses on the AD tasks, we note that there are some closely related problems, including out-of-distribution (OOD) detection [182] [183], novelty detection [116] [137], and open-set recognition (OSR) [51] [12]. Uniquely, AD usually does not assume the train set is anomaly-free, while other related tasks may do. Some methods designed for these related fields, e.g., OCSVM [157], can be used for AD as well; future benchmark can consider including: (i) OOD methods: MSP [65], energy-based EBO [104], and Mahalanobis distance-based MDS [92]; (ii) novelty detection methods: OCGAN [135] and Adversarial One-Class Classifier [154]; and (iii) OSR methods: OpenGAN [79] and PROSER [203]. See [155] for deeper connections and differences between AD and these fields.

We consider saliency detection (SD) [44] [46] and camouflage detection (CD) [45] as good inspirations and applications of AD tasks. Saliency detection identifies important regions in the images, where explainable AD algorithms [123], e.g., FCDD [106], may help the task. Camouflage detection finds concealed objects in the background, e.g., camouflaged anomalies blurred with normal objects [110], where camouflage-resistant AD methods [40] help detect concealed objects (that look normal but are abnormal). Future work can explore the explainability of detected objects in AD.

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 $\mathbf{O}_{\text{test}} := M(\mathbf{X}_{\text{test}})$.

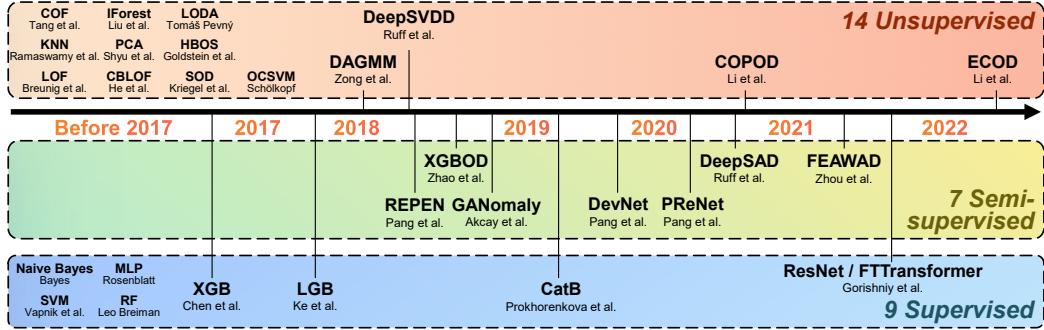


Figure 2: ADBench covers a wide range of AD algorithms. See Appx. B.1 for more details.

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 57 Datasets

Algorithms. Compared to the previous benchmarks, we have a larger algorithm collection with (i) the latest unsupervised AD algorithms like DeepSVDD [151] and ECOD [97]; (ii) SOTA semi-supervised algorithms, including DeepSAD [152] and DevNet [131]; (iii) latest network architectures like ResNet [62] in computer vision (CV) and Transformer [171] in the natural language processing (NLP) domain—we adapt ResNet and FTTransformer models [56] for tabular AD in the proposed ADBench; and (iv) ensemble learning methods like LightGBM [74], XGBoost [29], and CatBoost [138] that have shown effectiveness in AD tasks [170]. Fig. 2 shows the 30 algorithms (14 unsupervised, 7 semi-supervised, and 9 supervised algorithms) evaluated in ADBench, where we provide more information about them in Appx. B.1

Algorithm Implementation. Most unsupervised algorithms are readily available in our early work Python Outlier Detection (PyOD) [198], and some supervised methods are available in scikit-learn [133] and corresponding libraries. Supervised ResNet and FTTransformer tailored for tabular data have been open-sourced in their original paper [56]. We implement the semi-supervised methods and release them along with ADBench.

Public AD Datasets. In ADBench, we gather more than 40 benchmark datasets [25] [42] [129] [145], for model evaluation, as shown in Appx. Table B1. 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. For due diligence, we keep the datasets where the anomaly ratio is below 40% (Appx. Fig. B1).

Newly-added Datasets in ADBench. Since most of these datasets are relatively small, we introduce 10 more complex datasets from CV and NLP domains with more samples and richer features in ADBench (highlighted in Appx. Table B1). Pretrained models are applied to extract data embedding from CV and NLP datasets to access more complex representations, which has been widely used in AD literature [33] [115] [152] and shown better results than using the raw features. For NLP datasets, we use BERT [75] pretrained on the BookCorpus and English Wikipedia to extract the embedding of the [CLS] token. For CV datasets, we use ResNet18 [62] pretrained on the ImageNet [35] to extract the embedding after the last average pooling layer. Following previous works [151] [152], we set one of the multi-classes as normal, downsample the remaining classes to 5% of the total instances as anomalies, and report the average results over all the respective classes. Including these originally non-tabular datasets helps to see whether tabular AD methods can work on CV/NLP data after necessary preprocessing. See Appx. B.2 for more details on datasets.

3.3 Benchmark Angles in ADBench

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

Motivation. As shown in Table I existing benchmarks only focus on the unsupervised setting, i.e., none of the labeled anomalies is available. Despite, in addition to 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 [5] [7] [78] [189].

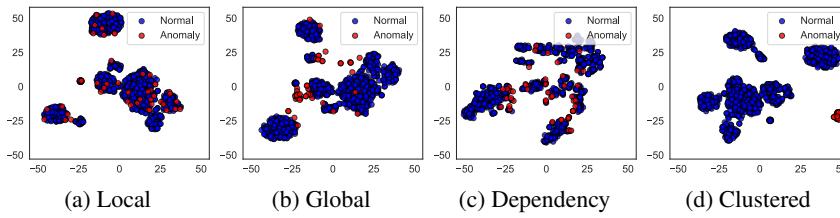


Figure 3: Illustration of four types of synthetic anomalies shown on Lymphography dataset. See the additional demo in Appx. Fig. B2

Notably, there is a group of semi-supervised AD algorithms [127] [128] [130] [131] [152] [168] [204] 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 levels of supervision following the settings in [127] [131] [204] to provide a fair comparison. For example, labeled anomalies $\gamma_l = 10\%$ means that 10% anomalies in the train set are known while other samples remain unlabeled. The complete 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 public 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 types of anomalies [55] [166]. In real-world applications, one may know specific types of anomalies of interest. To better understand the impact of anomaly types, we create synthetic datasets based on public datasets by injecting specific types of anomalies to analyze the response of AD algorithms.

Our design: In ADBench, we create *realistic* synthetic datasets from benchmark datasets by injecting specific types of anomalies. Some existing works, such as PyOD [198], generate fully synthetic anomalies by assuming their data distribution, which fails to create complex anomalies. We follow and enrich the approach in [166] 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 [166], Sparx [191], and ADBench) using the normal samples from a benchmark dataset and discard its original anomalies as we do not know their types. Then, We could generate normal samples and different types of anomalies based on their definitions by tweaking the generative model. The generation of normal samples is the same in all settings if not noted, and we provide the generation process of four types of anomalies below (also see our codebase for details).

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 [22]. We follow the GMM procedure [118] [166] 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 [68], 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 [117], i.e., the input features of dependency anomalies are assumed to be independent of each other. Vine Copula [1] method is applied to model the dependency structure of original data, where the probability density function of generated anomalies is set to complete independence by removing the modeled dependency (see [117]). We use Kernel Density Estimation (KDE) [61] to estimate the probability density function of features and generate normal samples.
- **Clustered anomalies**, also known as group anomalies [93], exhibit similar characteristics [42] [99]. 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 [169] 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 [24] [41] [47] [76] [177]. Meanwhile, the input data likely suffers from noise and corruption to some extent in real-world applications [42] [55] [60] [124]. 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 many applications, certain anomalies likely repeat multiple times in the data for reasons such as recording errors [83]. The presence of duplicated anomalies is also called the “anomaly masking” [55] [60] [100], posing challenges to many AD algorithms [25], e.g., the density-based KNN [11] [144]. Besides, the change of anomaly frequency would also affect the behavior of detection methods [42]. Therefore, we simulate this setting by splitting the data into train and test set, then duplicating the anomalies (both features and labels) up to 6 times in both sets, and observing how AD algorithms change.
- **Irrelevant Features.** Tabular data may contain irrelevant features caused by measurement noise or inconsistent measuring units [28] [55], where these noisy dimensions could hide the characteristics of anomaly data and thus make the detection process more difficult [128] [150]. 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 [131] [152] explored anomaly contamination in the unlabeled samples, we further discuss the more generalized impact of label contamination on the algorithm performance, where the label flips [122] [201] 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 98,436 experiments (Appx. C) to answer **Q1** (§4.2): How do AD algorithms perform 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 B1, ADBench includes 57 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 the newcomer data), we adapt all the algorithms for the inductive setting to predict the newcomer data, which is helpful in applications and also common in popular AD library PyOD [198], TODS [84], and PyGOD [102]. Thus, we use 70% data for training and the remaining 30% as the test 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 a 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) value.¹ Besides, the critical difference diagram (CD diagram) [34] [70] based on the Wilcoxon-Holm method is used for comparing groups of AD methods statistically ($p \leq 0.05$).

4.2 Overall Model Performance on Datasets with Varying Degrees of Supervision

As introduced in §3.3.1 we first present the results of unsupervised methods on 57 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.

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

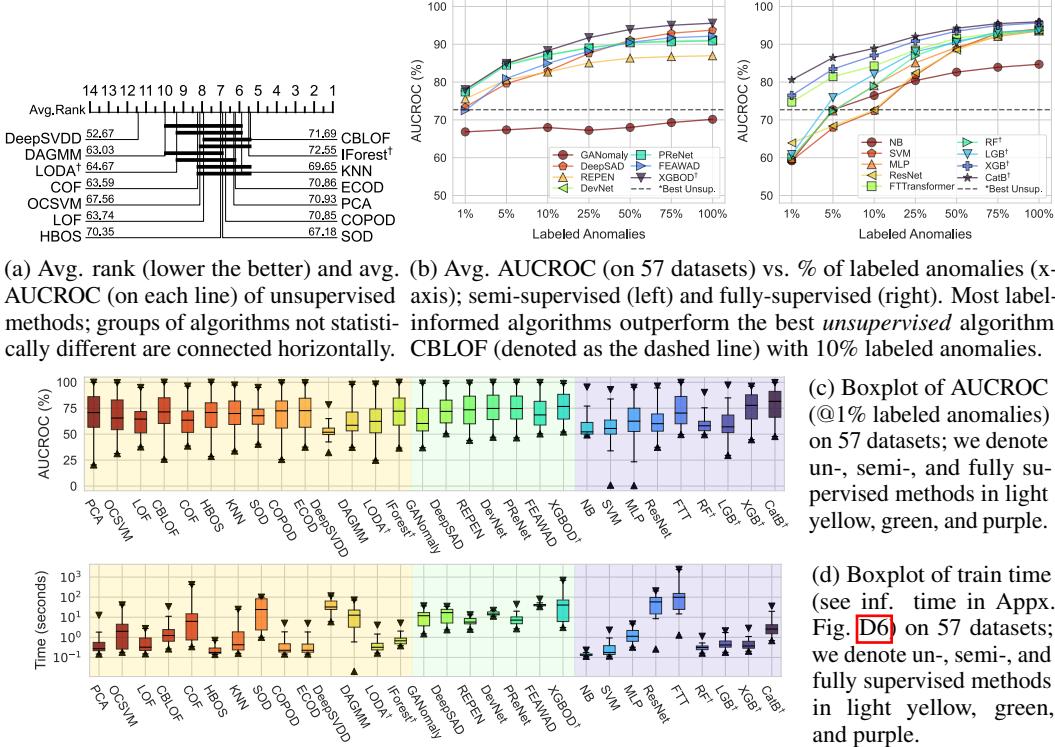


Figure 4: Average AD model performance across 57 benchmark datasets. (a) shows that no unsupervised algorithm statistically outperforms the rest. (b) shows that semi-supervised methods leverage the labels more efficiently than fully-supervised methods with a small labeled anomaly ratio γ_l . (c) and (d) present the boxplots of AUCROC and runtime. Ensemble methods are marked with "†".

None of the unsupervised methods is statistically better than the others, as shown in the critical difference diagram of Fig. 4a (where most algorithms are horizontally connected without statistical significance). We also note that some DL-based unsupervised methods like DeepSVDD and DAGMM are surprisingly worse than shallow methods. Without the guidance of label information, DL-based unsupervised algorithms are harder to train (due to more hyperparameters) and more difficult to tune 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= 75.56% for $\gamma_l = 1\%$ and AUCROC= 80.95% for $\gamma_l = 5\%$) are generally better than that of fully-supervised algorithms (median AUCROC= 60.84% for $\gamma_l = 1\%$ and AUCROC= 72.69% for $\gamma_l = 5\%$). For most 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 supervised methods need 10% labeled anomalies to achieve so. We also show the improvement of algorithm performances about the increasing γ_l , and notice that with a large number of labeled anomalies, both semi-supervised and supervised methods have comparable performance. Putting these together, we verify the assumed advantage of semi-supervised methods in leveraging limited label information more efficiently.

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 74.68% and 76.47%, respectively, outperforming the median AUCROC of all label-informed methods 72.91%. The great performance of tree-based ensembles (in tabular AD) is consistent with the findings in literature [20] [58] [170], which may be credited to their capacity to handle imbalanced AD datasets via aggregation. Future research may focus on understanding the cause and other merits of ensemble trees in tabular AD, e.g., better model efficiency.

Runtime Analysis. We present the train and inference time in Fig. 4d and Appx. Fig. D6. Runtime analysis finds that HBOS, COPOD, ECOD, and NB are the fastest as they treat each feature independently. In contrast, more complex representation learning methods like XGBOD, ResNet, and FITTransformer are computationally heavy. This should be factored in for algorithm selection.

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 difference diagrams). Meanwhile, the no-free-lunch theorem [175] suggests there is no universal winner for all tasks, and more focus should be spent on understanding the suitability of each AD algorithm. Notably, algorithm selection and hyperparameter optimization are important in unsupervised AD, but limited works [13] [109] [194] [199] have studied them. We may consider self-supervision [140] [158] [161] [179] and transfer learning [33] to improve tabular AD as well. Thus, we call for attention to 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 to semi-supervised AD methods which can leverage both the guidance from labels efficiently and the exploration of the unlabeled data. Regarding backbones, the 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.

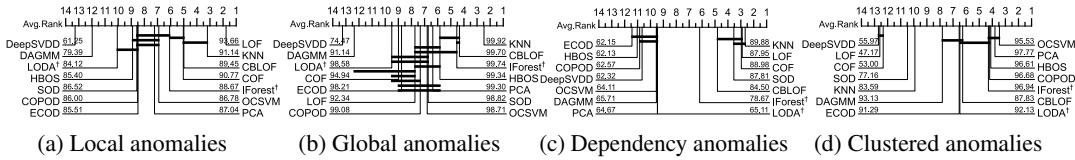


Figure 5: Avg. rank (lower the better) of unsupervised methods on different types of anomalies. Groups of algorithms not significantly different are connected horizontally in the CD diagrams. The unsupervised methods perform well when their assumptions conform to the underlying anomaly type.

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; LOF achieves the best AUCROC on local anomalies (Fig. 5a) and the second best AUCROC rank on dependency anomalies (Fig. 5c), but performs poorly on clustered anomalies (Fig. 5d). Practitioners should select algorithms based on the characteristics of the underlying task, and consider the algorithm which may cover more high-interest anomaly types [93].

The “power” of prior knowledge on anomaly types may outweigh the usage of partial labels. For the local, global, and dependency anomalies, most label-informed methods perform worse than the best unsupervised methods of each type (corresponding to LOF, KNN, and KNN). 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. 3a, 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 the 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 OCSVM, as few labeled anomalies can already represent 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 labels, and call for attention to designing anomaly-type-aware detection algorithms. In an

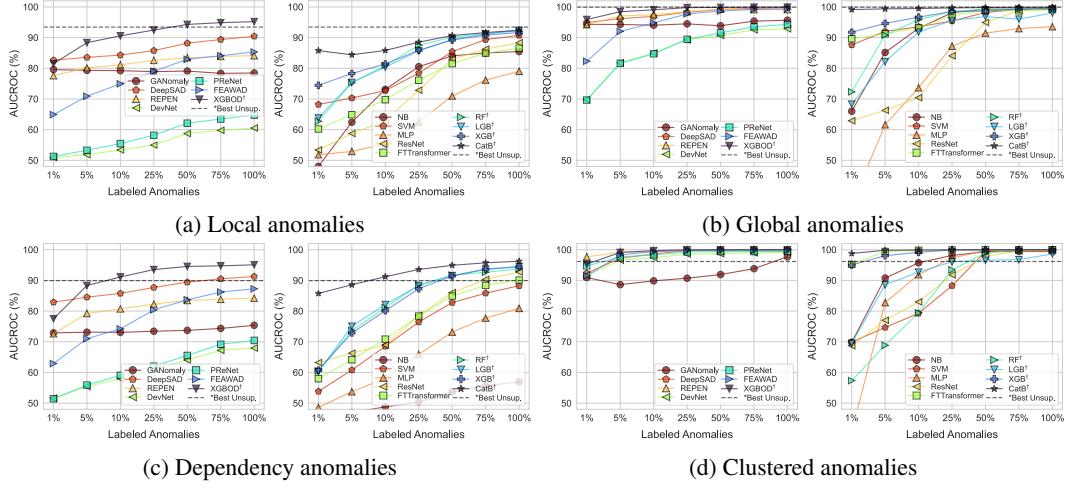


Figure 6: Semi- (left of each subfigure) and supervised (right) algorithms’ performance on different types of anomalies with varying levels of labeled anomalies. Surprisingly, these label-informed algorithms are *inferior* to the best unsupervised method except for the clustered anomalies.

ideal world, one may combine multiple AD algorithms based on the composition of anomaly types, via frameworks like dynamic model selection and combination [197]. To our knowledge, the latest advancement in this end [71] 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. Another interesting direction is to train an offline AD model using synthetically generated anomalies and then adapt it for online prediction on real-world datasets with likely similar anomaly types. Unsupervised domain adaption and transfer learning for AD [33] [185] may serve as useful references.

4.4 Algorithm Robustness under Noisy and Corrupted Data

In this section, we investigate the algorithm robustness (i.e., Δ performance; see absolute performance plot in Appx. D9) 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.

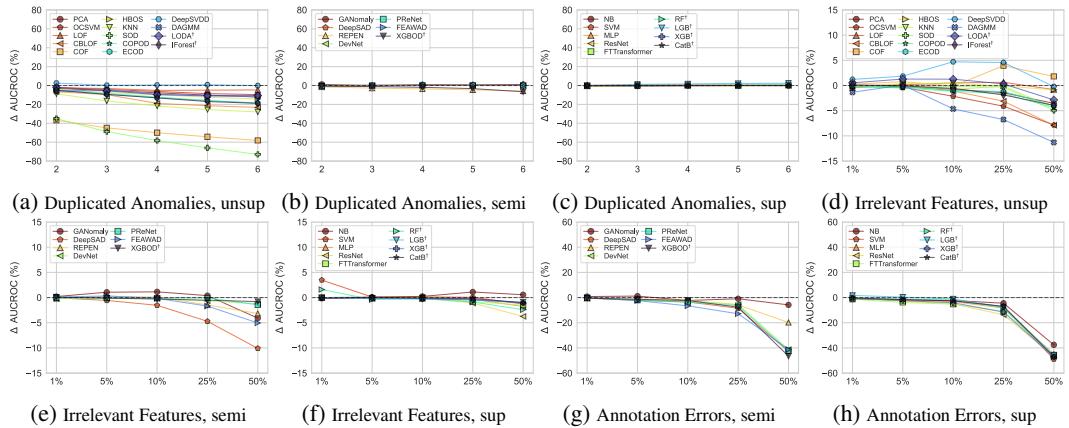


Figure 7: Algorithm performance change under noisy and corrupted data (i.e., duplicated anomalies for (a)-(c), irrelevant features for (d)-(f), and annotation errors for (g) and (h)). X-axis denotes either the duplicated times or the noise ratio. Y-axis denotes the % of performance change (Δ AUCROC), and its range remains consistent across different algorithms. The results reveal 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.

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 in duplication. When anomalies are duplicated by 6 times, the median Δ AUCROC of unsupervised methods is -16.43% , compared to that of semi-supervised methods -0.05% (Fig. 7b) and supervised methods 0.13% (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 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 performing 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 performing 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 that minor irrelevant features (e.g., 1%) help supervised methods as regularization to generalize better.

Both semi- and fully-supervised methods show great resilience to minor annotation errors. Although the detection performance of these methods is significantly downgraded 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 -1.52% and -1.91% , respectively. That being said, label-informed methods are still acceptable in practice as the annotation error should be relatively small [95][181].

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 [30][125][126] to combat irrelevant features. Moreover, label information could serve as effective guidance for 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 suggest using semi-supervised methods as the backbone for designing more robust AD algorithms. Also, recent works on leveraging multiple sets of noisy labels collectively for learning AD models are also relevant [200].

5 Conclusions and Future Work

In this paper, we introduce ADBench, the most comprehensive tabular anomaly detection benchmark with 30 algorithms and 57 benchmark datasets. Based on the analyses on multiple comparison angles, we unlock insights into the role of supervision, the importance of prior knowledge of 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 evaluating new algorithms.

ADBench can extend to understand the algorithm performance with (i) mixed types of anomalies; (ii) different levels of (intrinsic) anomaly ratio; and (iii) more data modalities. Also, future benchmarks can consider the latest algorithms [28][99][161], and curate datasets from emerging fields like drug discovery [69], molecule optimization [49][50], interpretability and explainability [123][180], and bias and fairness [32][67][123][159][165][190].

Acknowledgement

We briefly describe the authors' contributions. *Problem scoping:* M.J., Y.Z., S.H., X.H., and H.H.; *Experiment and Implementation:* M.J. and Y.Z.; *Result Analysis:* M.J., Y.Z., and X.H.; *Paper Drafting:* M.J., Y.Z., S.H., X.H., and H.H.; *Paper Revision:* M.J., Y.Z., S.H., and X.H.

We thank anonymous reviewers for their insightful feedback and comments. We appreciate the suggestions of Xueying Ding, Kwei-Herng (Henry) Lai, Meng-Chieh Lee, Ninghao Liu, Yuwen Yang, Allen Zhu, Chaochuan Hou, and Xu Yao. Y.Z. is partly supported by the Norton Graduate Fellowship. H.H., S.H., and M.J. are supported by the National Natural Science Foundation of China (NSFC) under Grant No. 72271151, 92146004.

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] S. Aeberhard, D. Coomans, and O. de Vel. The classification performance of rda. *Dept. of Computer Science and Dept. of Mathematics and Statistics, James Cook University of North Queensland, Tech. Rep*, pages 92–01, 1992.
- [3] C. C. Aggarwal. An introduction to outlier analysis. In *Outlier analysis*, pages 1–34. 2017.
- [4] C. C. Aggarwal. *Neural Networks and Deep Learning - A Textbook*. Springer, 2018.
- [5] N. B. Aissa and M. Guerroumi. Semi-supervised statistical approach for network anomaly detection. *Procedia Computer Science*, 83:1090–1095, 2016.
- [6] S. Akcay, D. Ameln, A. Vaidya, B. Lakshmanan, N. Ahuja, and U. Genc. Anomalib: A deep learning library for anomaly detection. *arXiv:2202.08341*, 2022.
- [7] S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. In *ACCV*, pages 622–637, 2018.
- [8] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon. Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection. In *IJCNN*, pages 1–8. IEEE, 2019.
- [9] F. Alimoglu and E. Alpaydin. Methods of combining multiple classifiers based on different representations for pen-based handwritten digit recognition. In *TAINN*. Citeseer, 1996.
- [10] E. Alpaydin and C. Kaynak. Cascading classifiers. *Kybernetika*, 34(4):369–374, 1998.
- [11] F. Angiulli and C. Pizzuti. Fast outlier detection in high dimensional spaces. In *ECML/PKDD*, pages 15–27. Springer, 2002.
- [12] D. Ayres-de Campos, J. Bernardes, A. Garrido, J. Marques-de Sa, and L. Pereira-Leite. Sisporto 2.0: a program for automated analysis of cardiotocograms. *Journal of Maternal-Fetal Medicine*, 2000.
- [13] M. Bahri, F. Salutari, A. Putina, and M. Sozio. Automl: state of the art with a focus on anomaly detection, challenges, and research directions. *IDSA*, pages 1–14, 2022.
- [14] T. Bayes. Lii. an essay towards solving a problem in the doctrine of chances. *Philos. Trans. Royal Soc. A*, pages 370–418, 1763.
- [15] L. Bergman and Y. Hoshen. Classification-based anomaly detection for general data. In *ICLR*, 2019.
- [16] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger. Mttec ad—a comprehensive real-world dataset for unsupervised anomaly detection. In *CVPR*, pages 9592–9600, 2019.
- [17] E. Berthonnaud, J. Dimnet, P. Roussouly, and H. Labelle. Analysis of the sagittal balance of the spine and pelvis using shape and orientation parameters. *Clinical Spine Surgery*, 18(1):40–47, 2005.
- [18] J. A. Blackard and D. J. Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and electronics in agriculture*, 24(3):131–151, 1999.
- [19] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov. Enriching word vectors with subword information. *TACL*, 5:135–146, 2017.
- [20] V. Borisov, T. Leemann, K. Seßler, J. Haug, M. Pawelczyk, and G. Kasneci. Deep neural networks and tabular data: A survey. *arXiv:2110.01889*, 2021.
- [21] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [22] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. Lof: identifying density-based local outliers. In *SIGMOD*, pages 93–104, 2000.
- [23] N. Brümmer, S. Cumaní, O. Glembek, M. Karafiát, P. Matějka, J. Pešán, O. Plchot, M. Soufifar, E. d. Villiers, and J. H. Černocký. Description and analysis of the brno276 system for lre2011. In *Odyssey 2012-the speaker and language recognition workshop*, 2012.
- [24] H. Cai, J. Liu, and W. Yin. Learned robust pca: A scalable deep unfolding approach for high-dimensional outlier detection. *NeurIPS*, 34, 2021.
- [25] 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.
- [26] B. Cestnik, I. Kononenko, and I. Bratko. A knowledge-elicitation tool for sophisticated users. In *European Conference on European Working Session on Learning EWSL*, volume 87, 1987.
- [27] 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 *NeurIPS*, 2021.

- [28] C.-H. Chang, J. Yoon, S. Arik, M. Udell, and T. Pfister. Data-efficient and interpretable tabular anomaly detection. *ArXiv*, 2203.02034, 2022.
- [29] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *KDD*, pages 785–794, 2016.
- [30] L. Cheng, Y. Wang, X. Liu, and B. Li. Outlier detection ensemble with embedded feature selection. In *AAAI*, volume 34, pages 3503–3512, 2020.
- [31] C. Cortes and V. Vapnik. Support vector machine. *Machine learning*, 20(3):273–297, 1995.
- [32] I. Davidson and S. S. Ravi. A framework for determining the fairness of outlier detection. In *ECAI 2020*, pages 2465–2472. IOS Press, 2020.
- [33] L. Deecke, L. Ruff, R. A. Vandermeulen, and H. Bilen. Transfer-based semantic anomaly detection. In *ICML*, pages 2546–2558, 2021.
- [34] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *The JMLR*, 7:1–30, 2006.
- [35] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009.
- [36] P. Diaconis and B. Efron. Computer-intensive methods in statistics. *Scientific American*, 248(5), 1983.
- [37] T. Dietterich, A. Jain, R. Lathrop, and T. Lozano-Perez. A comparison of dynamic reposing and tangent distance for drug activity prediction. *NeurIPS*, 6, 1993.
- [38] 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.
- [39] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.
- [40] Y. Dou, Z. Liu, L. Sun, Y. Deng, H. Peng, and P. S. Yu. Enhancing graph neural network-based fraud detectors against camouflaged fraudsters. In *CIKM*, pages 315–324, 2020.
- [41] 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 *ICML*, pages 2880–2891, 2021.
- [42] A. Emmott, S. Das, T. Dietterich, A. Fern, and W.-K. Wong. A meta-analysis of the anomaly detection problem. *ArXiv*, 1503.01158, 2015.
- [43] I. W. Evett and E. J. Spiehler. Rule induction in forensic science. pages 152–160, 1989.
- [44] D.-P. Fan, M.-M. Cheng, J.-J. Liu, S.-H. Gao, Q. Hou, and A. Borji. Salient objects in clutter: Bringing salient object detection to the foreground. In *ECCV*, pages 186–202, 2018.
- [45] D.-P. Fan, G.-P. Ji, G. Sun, M.-M. Cheng, J. Shen, and L. Shao. Camouflaged object detection. In *CVPR*, pages 2777–2787, 2020.
- [46] D.-P. Fan, T. Li, Z. Lin, G.-P. Ji, D. Zhang, M.-M. Cheng, H. Fu, and J. Shen. Re-thinking co-salient object detection. *TPAMI*, 2021.
- [47] L. Fan, S. Liu, P.-Y. Chen, G. Zhang, and C. Gan. When does contrastive learning preserve adversarial robustness from pretraining to finetuning? *NeurIPS*, 34, 2021.
- [48] P. W. Frey and D. J. Slate. Letter recognition using holland-style adaptive classifiers. *Machine learning*, 6(2):161–182, 1991.
- [49] T. Fu, C. Xiao, X. Li, L. M. Glass, and J. Sun. Mimosa: Multi-constraint molecule sampling for molecule optimization. In *AAAI*, volume 35, pages 125–133, 2021.
- [50] T. Fu, C. Xiao, and J. Sun. Core: Automatic molecule optimization using copy & refine strategy. In *AAAI*, volume 34, pages 638–645, 2020.
- [51] C. Geng, S.-j. Huang, and S. Chen. Recent advances in open set recognition: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3614–3631, 2020.
- [52] M. Goldstein and A. Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: poster and demo track*, 9, 2012.
- [53] M. Goldstein and S. Uchida. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *Plos one*, 11(4):e0152173, 2016.
- [54] I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016.
- [55] P. Gopalan, V. Sharan, and U. Wieder. Pidforest: anomaly detection via partial identification. *NeurIPS*, 32, 2019.
- [56] Y. Gorishniy, I. Rubachev, V. Khrulkov, and A. Babenko. Revisiting deep learning models for tabular data. *NeurIPS*, 34, 2021.

- [57] N. Görnitz, M. Kloft, K. Rieck, and U. Brefeld. Toward supervised anomaly detection. *JAIR*, 46:235–262, 2013.
- [58] L. Grinsztajn, E. Oyallon, and G. Varoquaux. Why do tree-based models still outperform deep learning on tabular data? *arXiv:2207.08815*, 2022.
- [59] C. Grunau and V. Rozhoň. Adapting k-means algorithms for outliers. *ArXiv*, 2007.01118, 2020.
- [60] S. Guha, N. Mishra, G. Roy, and O. Schrijvers. Robust random cut forest based anomaly detection on streams. In *ICML*, pages 2712–2721, 2016.
- [61] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. 2009.
- [62] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [63] R. He and J. McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *WWW*, pages 507–517, 2016.
- [64] Z. He, X. Xu, and S. Deng. Discovering cluster-based local outliers. *Pattern recognition letters*, 24(9–10):1641–1650, 2003.
- [65] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *ICLR*, 2017.
- [66] P. Horton and K. Nakai. A probabilistic classification system for predicting the cellular localization sites of proteins. In *Ismb*, volume 4, pages 109–115, 1996.
- [67] X. Hu, Y. Huang, B. Li, and T. Lu. Uncovering the source of machine bias. *KDD 2021, Machine Learning for Consumers and Markets Workshop*, 2021.
- [68] H. Huang, H. Qin, S. Yoo, and D. Yu. Physics-based anomaly detection defined on manifold space. *TKDD*, 9(2):1–39, 2014.
- [69] 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. *NeurIPS*, 2021.
- [70] 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.
- [71] C. I. Jerez, J. Zhang, and M. R. Silva. On equivalence of anomaly detection algorithms. *TKDD*, 2022.
- [72] A. Joulin, É. Grave, P. Bojanowski, and T. Mikolov. Bag of tricks for efficient text classification. In *EACL*, pages 427–431, 2017.
- [73] Y. Kawachi, Y. Koizumi, and N. Harada. Complementary set variational autoencoder for supervised anomaly detection. In *ICASSP*, pages 2366–2370. IEEE, 2018.
- [74] 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. *NeurIPS*, 30:3146–3154, 2017.
- [75] 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.
- [76] M. Kim, J. Tack, and S. J. Hwang. Adversarial self-supervised contrastive learning. *NeurIPS*, 33:2983–2994, 2020.
- [77] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *ArXiv*, 1412.6980, 2014.
- [78] 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.
- [79] S. Kong and D. Ramanan. Opengan: Open-set recognition via open data generation. In *CVPR*, pages 813–822, 2021.
- [80] H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. Outlier detection in axis-parallel subspaces of high dimensional data. In *PAKDD*, pages 831–838. Springer, 2009.
- [81] A. Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, University of Tront*, 2009.
- [82] M. Kudo, J. Toyama, and M. Shimbo. Multidimensional curve classification using passing-through regions. *Pattern Recognition Letters*, 20(11-13):1103–1111, 1999.
- [83] D. Kwon, K. Natarajan, S. C. Suh, H. Kim, and J. Kim. An empirical study on network anomaly detection using convolutional neural networks. In *ICDCS*, pages 1595–1598, 2018.
- [84] 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 *AAAI*, pages 16060–16062, 2021.

- [85] K.-H. Lai, D. Zha, J. Xu, Y. Zhao, G. Wang, and X. Hu. Revisiting time series outlier detection: Definitions and benchmarks. In *NeurIPS*, 2021.
- [86] K. Lang. Newsweeder: Learning to filter netnews. In *ICML*, pages 331–339. Elsevier, 1995.
- [87] A. Lavin and S. Ahmad. Evaluating real-time anomaly detection algorithms—the numenata anomaly benchmark. In *2015 IEEE 14th ICML and applications (ICMLA)*, pages 38–44. IEEE, 2015.
- [88] A. Lazarevic, L. Ertoz, V. Kumar, A. Ozgur, and J. Srivastava. A comparative study of anomaly detection schemes in network intrusion detection. In *SDM*, pages 25–36. SIAM, 2003.
- [89] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [90] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [91] C. H. Lee and K. Lee. Semi-supervised anomaly detection algorithm based on kl divergence (sad-kl). *arXiv:2203.14539*, 2022.
- [92] K. Lee, K. Lee, H. Lee, and J. Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *NeurIPS*, 31, 2018.
- [93] M.-C. Lee, S. Shekhar, C. Faloutsos, T. N. Hutson, and L. Iasemidis. Gen 2 out: Detecting and ranking generalized anomalies. In *Big Data*, pages 801–811. IEEE, 2021.
- [94] 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 *Big Data*, pages 950–956. IEEE, 2020.
- [95] G. Li, Y. Xie, and L. Lin. Weakly supervised salient object detection using image labels. In *AAAI*, volume 32, 2018.
- [96] Z. Li, Y. Zhao, N. Botta, C. Ionescu, and X. Hu. Copod: copula-based outlier detection. In *ICDM*, pages 1118–1123. IEEE, 2020.
- [97] Z. Li, Y. Zhao, X. Hu, N. Botta, C. Ionescu, and G. Chen. Ecod: Unsupervised outlier detection using empirical cumulative distribution functions. *TKDE*, pages 1–1, 2022.
- [98] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017.
- [99] B. Liu, P. Tan, and J. Zhou. Unsupervised anomaly detection by robust density estimation. In *AAAI*, pages 4101–4108. AAAI Press, 2022.
- [100] F. T. Liu, K. M. Ting, and Z.-H. Zhou. Isolation forest. In *ICDM*, pages 413–422. IEEE, 2008.
- [101] K. Liu, Y. Dou, Y. Zhao, X. Ding, X. Hu, R. Zhang, K. Ding, C. Chen, H. Peng, K. Shu, et al. Benchmarking node outlier detection on graphs. *arXiv preprint arXiv:2206.10071*, 2022.
- [102] 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*, 2204.12095, 2022.
- [103] S. Liu and M. Hauskrecht. Event outlier detection in continuous time. In *ICML*, pages 6793–6803, 2021.
- [104] W. Liu, X. Wang, J. Owens, and Y. Li. Energy-based out-of-distribution detection. *NeurIPS*, 33:21464–21475, 2020.
- [105] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv:1907.11692*, 2019.
- [106] P. Liznerski, L. Ruff, R. A. Vandermeulen, B. J. Franks, M. Kloft, and K. R. Muller. Explainable deep one-class classification. In *ICLR*, 2020.
- [107] W.-Y. Loh. Classification and regression trees. *WIREs Data Mining and Knowledge Discovery*, 1, 2011.
- [108] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *ArXiv*, 1711.05101, 2017.
- [109] M. Q. Ma, Y. Zhao, X. Zhang, and L. Akoglu. A large-scale study on unsupervised outlier model selection: Do internal strategies suffice? *ArXiv*, 2104.01422, 2021.
- [110] X. Ma, J. Wu, S. Xue, J. Yang, C. Zhou, Q. Z. Sheng, H. Xiong, and L. Akoglu. A comprehensive survey on graph anomaly detection with deep learning. *TKDE*, 2021.
- [111] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts. Learning word vectors for sentiment analysis. In *ACL*, pages 142–150. Association for Computational Linguistics, 2011.
- [112] A. Mahdavi and M. Carvalho. A survey on open set recognition. In *AIKE*, pages 37–44. IEEE, 2021.
- [113] D. Malerba, F. Esposito, and G. Semeraro. A further comparison of simplification methods for decision-tree induction. In *Learning from data*, pages 365–374. Springer, 1996.
- [114] O. L. Mangasarian, W. N. Street, and W. H. Wolberg. Breast cancer diagnosis and prognosis via linear programming. *Operations Research*, 43(4):570–577, 1995.

- [115] A. Manolache, F. Brad, and E. Burceanu. Date: Detecting anomalies in text via self-supervision of transformers. In *NAACL*, pages 267–277, 2021.
- [116] M. Markou and S. Singh. Novelty detection: a review—part 1: statistical approaches. *Signal processing*, 83(12):2481–2497, 2003.
- [117] R. Martinez-Guerra and J. L. Mata-Machuca. *Fault detection and diagnosis in nonlinear systems*. Springer, 2016.
- [118] G. W. Milligan. An algorithm for generating artificial test clusters. *Psychometrika*, 50(1):123–127, 1985.
- [119] N. Moustafa and J. Slay. Unsw-nb15: a comprehensive data set for network intrusion detection systems (unsw-nb15 network data set). In *2015 military communications and information systems conference (MilCIS)*, pages 1–6. IEEE, 2015.
- [120] N. Mu and J. Gilmer. Mnist-c: A robustness benchmark for computer vision. *arXiv:1906.02337*, 2019.
- [121] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [122] D. T. Nguyen, C. K. Mummadi, T. P. N. Ngo, T. H. P. Nguyen, L. Begel, and T. Brox. Self: Learning to filter noisy labels with self-ensembling. In *ICLR*, 2019.
- [123] G. Pang and C. Aggarwal. Toward explainable deep anomaly detection. In *KDD*, pages 4056–4057, 2021.
- [124] 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.
- [125] G. Pang, L. Cao, L. Chen, and H. Liu. Unsupervised feature selection for outlier detection by modelling hierarchical value-feature couplings. In *ICDM*, pages 410–419. IEEE, 2016.
- [126] 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.
- [127] G. Pang, L. Cao, L. Chen, and H. Liu. Learning representations of ultrahigh-dimensional data for random distance-based outlier detection. In *KDD*, pages 2041–2050, 2018.
- [128] G. Pang, C. Ding, C. Shen, and A. v. d. Hengel. Explainable deep few-shot anomaly detection with deviation networks. *ArXiv*, 2108.00462, 2021.
- [129] 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.
- [130] G. Pang, C. Shen, H. Jin, and A. v. d. Hengel. Deep weakly-supervised anomaly detection. *ArXiv*, 1910.13601, 2019.
- [131] G. Pang, C. Shen, and A. van den Hengel. Deep anomaly detection with deviation networks. In *KDD*, pages 353–362, 2019.
- [132] J. Paparrizos, Y. Kang, P. Boniol, R. S. Tsay, T. Palpanas, and M. J. Franklin. Tsb-uad: an end-to-end benchmark suite for univariate time-series anomaly detection. *VLDB*, 15(8):1697–1711, 2022.
- [133] 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 JMLR*, 12:2825–2830, 2011.
- [134] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *EMNLP*, pages 1532–1543, 2014.
- [135] P. Perera, R. Nallapati, and B. Xiang. Ogan: One-class novelty detection using gans with constrained latent representations. In *CVPR*, pages 2898–2906, 2019.
- [136] T. Pevný. Loda: Lightweight on-line detector of anomalies. *Machine Learning*, 102(2):275–304, 2016.
- [137] M. A. Pimentel, D. A. Clifton, L. Clifton, and L. Tarassenko. A review of novelty detection. *Signal processing*, 99:215–249, 2014.
- [138] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, and A. Gulin. Catboost: unbiased boosting with categorical features. *NeurIPS*, 31, 2018.
- [139] C. Qiu, A. Li, M. Kloft, M. Rudolph, and S. Mandt. Latent outlier exposure for anomaly detection with contaminated data. *ArXiv*, 2202.08088, 2022.
- [140] C. Qiu, T. Pfrommer, M. Kloft, S. Mandt, and M. Rudolph. Neural transformation learning for deep anomaly detection beyond images. In *ICML*, pages 8703–8714, 2021.
- [141] J. R. Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.
- [142] J. R. Quinlan, P. J. Compton, K. Horn, and L. Lazarus. Inductive knowledge acquisition: a case study. In *Australian Conference on Applications of expert systems*, pages 137–156, 1987.
- [143] M. M. Rahman, D. Balakrishnan, D. Murthy, M. Kutlu, and M. Lease. An information retrieval approach to building datasets for hate speech detection. In *NeurIPS*, 2021.

- [144] S. Ramaswamy, R. Rastogi, and K. Shim. Efficient algorithms for mining outliers from large data sets. In *SIGMOD*, pages 427–438, 2000.
- [145] S. Rayana. ODDS library, 2016.
- [146] Q. Rebjock, B. Kurt, T. Januschowski, and L. Callot. Online false discovery rate control for anomaly detection in time series. *NeurIPS*, 34, 2021.
- [147] T. Reiss, N. Cohen, L. Bergman, and Y. Hoshen. Panda: Adapting pretrained features for anomaly detection and segmentation. In *CVPR*, pages 2806–2814, 2021.
- [148] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- [149] S. Ruder. An overview of gradient descent optimization algorithms. *ArXiv*, 1609.04747, 2016.
- [150] 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.
- [151] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft. Deep one-class classification. In *ICML*, pages 4393–4402, 2018.
- [152] L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K. Müller, and M. Kloft. Deep semi-supervised anomaly detection. In *ICLR*. OpenReview.net, 2020.
- [153] L. Ruff, Y. Zemlyanskiy, R. Vandermeulen, T. Schnake, and M. Kloft. Self-attentive, multi-context one-class classification for unsupervised anomaly detection on text. In *ACL*, pages 4061–4071, 2019.
- [154] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli. Adversarially learned one-class classifier for novelty detection. In *CVPR*, pages 3379–3388, 2018.
- [155] M. Salehi, H. Mirzaei, D. Hendrycks, Y. Li, M. H. Rohban, and M. Sabokrou. A unified survey on anomaly, novelty, open-set, and out-of-distribution detection: Solutions and future challenges. *arXiv:2110.14051*, 2021.
- [156] R. Schirrmeister, Y. Zhou, T. Ball, and D. Zhang. Understanding anomaly detection with deep invertible networks through hierarchies of distributions and features. *NeurIPS*, 33:21038–21049, 2020.
- [157] 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.
- [158] V. Sehwag, M. Chiang, and P. Mittal. Ssd: A unified framework for self-supervised outlier detection. In *ICLR*, 2020.
- [159] S. Shekhar, N. Shah, and L. Akoglu. Fairod: Fairness-aware outlier detection. In *AIES*, pages 210–220, 2021.
- [160] L. Shen, Z. Li, and J. Kwok. Timeseries anomaly detection using temporal hierarchical one-class network. *NeurIPS*, 33:13016–13026, 2020.
- [161] T. Shenkar and L. Wolf. Anomaly detection for tabular data with internal contrastive learning. In *International Conference on Learning Representations*, 2021.
- [162] 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 FL Dept of Electrical and Computer Engineering, 2003.
- [163] V. G. Sigillito, S. P. Wing, L. V. Hutton, and K. B. Baker. Classification of radar returns from the ionosphere using neural networks. *Johns Hopkins APL Technical Digest*, 10(3):262–266, 1989.
- [164] J. Soenen, E. Van Wolputte, L. Perini, V. Verheyen, 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, 2021.
- [165] H. Song, P. Li, and H. Liu. Deep clustering based fair outlier detection. In *KDD*, pages 1481–1489, 2021.
- [166] G. Steinbuss and K. Böhm. Benchmarking unsupervised outlier detection with realistic synthetic data. *TKDD*, 15(4):1–20, 2021.
- [167] J. Tang, Z. Chen, A. W.-C. Fu, and D. W. Cheung. Enhancing effectiveness of outlier detections for low density patterns. In *PAKDD*, pages 535–548. Springer, 2002.
- [168] B. Tian, Q. Su, and J. Yin. Anomaly detection by leveraging incomplete anomalous knowledge with anomaly-aware bidirectional gans. *ArXiv*, 2204.13335, 2022.
- [169] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *JMLR*, 9(11), 2008.
- [170] S. Vargaftik, I. Keslassy, A. Orda, and Y. Ben-Itzhak. Rade: Resource-efficient supervised anomaly detection using decision tree-based ensemble methods. *Machine Learning*, 110(10):2835–2866, 2021.
- [171] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017.

- [172] Z. Wang, B. Dai, D. Wipf, and J. Zhu. Further analysis of outlier detection with deep generative models. *NeurIPS*, 33:8982–8992, 2020.
- [173] W. H. Wolberg and O. L. Mangasarian. Multisurface method of pattern separation for medical diagnosis applied to breast cytology. *Proceedings of the national academy of sciences*, 87(23):9193–9196, 1990.
- [174] D. H. Wolpert. Stacked generalization. *Neural networks*, 5(2):241–259, 1992.
- [175] D. H. Wolpert and W. G. Macready. No free lunch theorems for optimization. *IEEE transactions on evolutionary computation*, 1(1):67–82, 1997.
- [176] K. S. Woods, J. L. Solka, C. E. Priebe, W. P. Kegelmeyer Jr, C. C. Doss, and K. W. Bowyer. Comparative evaluation of pattern recognition techniques for detection of microcalcifications in mammography. In *State of The Art in Digital Mammographic Image Analysis*, pages 213–231. World Scientific, 1994.
- [177] B. Wu, J. Chen, D. Cai, X. He, and Q. Gu. Do wider neural networks really help adversarial robustness? *NeurIPS*, 34, 2021.
- [178] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv:1708.07747*, 2017.
- [179] Z. Xiao, Q. Yan, and Y. Amit. Do we really need to learn representations from in-domain data for outlier detection? *ArXiv*, 2105.09270, 2021.
- [180] H. Xu, Y. Wang, S. Jian, Z. Huang, Y. Wang, N. Liu, and F. Li. Beyond outlier detection: Outlier interpretation by attention-guided triplet deviation network. In *WWW*, pages 1328–1339, 2021.
- [181] Y. Xu, J. Ding, L. Zhang, and S. Zhou. Dp-ssl: Towards robust semi-supervised learning with a few labeled samples. *NeurIPS*, 34, 2021.
- [182] J. Yang, P. Wang, D. Zou, Z. Zhou, K. Ding, W. Peng, H. Wang, G. Chen, B. Li, Y. Sun, et al. Openood: Benchmarking generalized out-of-distribution detection. 2022.
- [183] J. Yang, K. Zhou, Y. Li, and Z. Liu. Generalized out-of-distribution detection: A survey. *arXiv:2110.11334*, 2021.
- [184] L. Yang, Z. Zhang, S. Hong, R. Xu, Y. Zhao, Y. Shao, W. Zhang, M.-H. Yang, and B. Cui. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022.
- [185] Z. Yang, I. S. Bozchalooi, and E. Darve. Anomaly detection with domain adaptation. *arXiv:2006.03689*, 2020.
- [186] 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.
- [187] M. D. Zeiler. Adadelta: an adaptive learning rate method. *ArXiv*, 1212.5701, 2012.
- [188] H. Zenati, M. Romain, C.-S. Foo, B. Lecouat, and V. Chandrasekhar. Adversarially learned anomaly detection. In *ICDM*, pages 727–736. IEEE, 2018.
- [189] D. Zha, K.-H. Lai, M. Wan, and X. Hu. Meta-aad: Active anomaly detection with deep reinforcement learning. In *ICDM*, pages 771–780. IEEE, 2020.
- [190] H. Zhang and I. Davidson. Towards fair deep anomaly detection. In *FAccT*, pages 138–148, 2021.
- [191] S. Zhang, V. Ursekar, and L. Akoglu. Sparx: Distributed outlier detection at scale. *KDD*, 2022.
- [192] X. Zhang, J. Zhao, and Y. LeCun. Character-level convolutional networks for text classification. *NIPS*, 28, 2015.
- [193] 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.
- [194] Y. Zhao and L. Akoglu. Towards unsupervised hpo for outlier detection. *arXiv preprint arXiv:2208.11727*, 2022.
- [195] Y. Zhao and M. K. Hryniwicki. Xgbod: improving supervised outlier detection with unsupervised representation learning. In *IJCNN*, pages 1–8. IEEE, 2018.
- [196] 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. *MLSys*, 3:463–478, 2021.
- [197] Y. Zhao, Z. Nasrullah, M. K. Hryniwicki, and Z. Li. Lscp: Locally selective combination in parallel outlier ensembles. In *SDM*, pages 585–593. SIAM, 2019.
- [198] Y. Zhao, Z. Nasrullah, and Z. Li. Pyod: A python toolbox for scalable outlier detection. *JMLR*, 20:1–7, 2019.
- [199] Y. Zhao, R. Rossi, and L. Akoglu. Automatic unsupervised outlier model selection. *NeurIPS*, 34, 2021.
- [200] Y. Zhao, G. Zheng, S. Mukherjee, R. McCann, and A. Awadallah. Admoe: Anomaly detection with mixture-of-experts from noisy labels. *arXiv preprint arXiv:2208.11290*, 2022.

- [201] G. Zheng, A. H. Awadallah, and S. Dumais. Meta label correction for noisy label learning. *AAAI*, 2021.
- [202] Y. Zheng, X. Wang, Y. Qi, W. Li, and L. Wu. Benchmarking unsupervised anomaly detection and localization. *ArXiv*, 2205.14852, 2022.
- [203] D.-W. Zhou, H.-J. Ye, and D.-C. Zhan. Learning placeholders for open-set recognition. In *CVPR*, pages 4401–4410, 2021.
- [204] Y. Zhou, X. Song, Y. Zhang, F. Liu, C. Zhu, and L. Liu. Feature encoding with autoencoders for weakly supervised anomaly detection. *TNNLS*, 2021.
- [205] Z.-H. Zhou. A brief introduction to weakly supervised learning. *Natl. Sci. Rev.*, 5(1):44–53, 2018.
- [206] 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 *ICLR*, 2018.

Supplementary Material for ADBench: Anomaly Detection Benchmark

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

A Related Works with More Details

We provide more details on existing AD algorithms and benchmarks, and 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 [3], e.g., anomalies located 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 [3] [15] [129] [150] [198], which can be roughly categorized into shallow and deep (neural network) methods. The former often carry better interpretability, while the latter handles large, high-dimensional data better. Recent book [3] and surveys [129] [150] 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.1.

Representative Shallow Methods. Some representative shallow methods include: (i) Isolation Forest (IForest) [100] 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) [157] 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) [97] 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 have gained more attention recently, and we briefly review some representative ones in this section. Deep Autoencoding Gaussian Mixture Model (DAGMM) [206] 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 pre-training. Deep Support Vector Data Description (DeepSVDD) [151] 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 [129]. Although some loss functions (e.g., focal loss [98]) are devised to address the class imbalance problem, they are often not specific for AD tasks. There also exist a few works [57] [73] discussing the relationship between fully-supervised and semi-supervised AD methods, and argue that semi-supervised AD needs to be 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.1), and recommend interesting readers to recent machine learning books [4] [54] and scikit-learn [133] 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 to detect unseen types of anomalies. To this end, some recent studies investigate efficiently using 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.1 for more algorithm details and hyperparameter settings in ADBench.

Representative Methods. Extreme Gradient Boosting Outlier Detection (XGBOD) [195] 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) [152] is an end-to-end methodology considered the state-of-the-art method in semi-supervised anomaly detection. DeepSAD improves the DeepSVDD [151] model by the inverse loss function for the labeled anomalies. REPresentations for a random nEarest Neighbor distance-based method (REPEN) [127] 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) [131] constructs an end-to-end neural network for learning anomaly scores, which forces the network to produce statistically higher anomaly scores for identified anomalies than that of unlabeled data. Pairwise Relation prediction-based ordinal regression Network (PReNet) [130] 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) [204] 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 [131] and DAGMM [206].

A.4 Existing AD Benchmarks

As we show in Table 1, there is a line of existing AD benchmarks. [150] discusses a unifying review of both the shallow and deep anomaly detection methods, but they mainly focus on the theoretical perspective and thus lack results from the experimental views. [25] benchmarks 19 different unsupervised methods on 10 datasets, and analyzes the characteristics of density-based and clustering-based algorithms. [38] tests 14 unsupervised anomaly detection methods on 15 public datasets, and analyzes the scalability, memory consumption, and robustness of different methods. [166] 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. [42] evaluates 8 unsupervised methods on 19 public datasets, and produces a large corpus of synthetic anomaly detection datasets that vary in their construction across several dimensions that are important to real-world applications. [25] compares 12 unsupervised anomaly detection approaches on 23 datasets, providing a 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 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)** [162]. 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).

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

²<https://scikit-learn.org/stable/>

2. **One-class SVM (OCSVM)** [157]. OCSVM maximizes the margin between the origin and the normal samples, and defines the decision boundary as the hyperplane that determines the margin.
3. **Local Outlier Factor (LOF)** [22]. LOF measures the local deviation of the density of a given sample with respect to its neighbors.
4. **Clustering Based Local Outlier Factor (CBLOF)** [64]. 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)** [167]. COF uses the ratio of the average chaining distance of data points and the average chaining distance of k -th nearest neighbor of the data point, as the anomaly score for observations.
6. **Histogram-based outlier detection (HBOS)** [52]. HBOS assumes feature independence and calculates the degree of outlyingness by building histograms.
7. **K-Nearest Neighbors (KNN)** [144]. KNN views the anomaly score of the input instance as the distance to its k -th nearest neighbor.
8. **Subspace Outlier Detection (SOD)** [80]. SOD aims to detect outliers in varying subspaces of high-dimensional feature space.
9. **Copula Based Outlier Detector (COPOD)** [96]. COPOD is a hyperparameter-free, highly interpretable outlier detection algorithm based on empirical copula models.
10. **Empirical-Cumulative-distribution-based Outlier Detection (ECOD)** [97]. ECOD is a hyperparameter-free, highly interpretable outlier detection algorithm based on empirical CDF functions. Basically, it uses ECDF to estimate the density of each feature independently, and assumes that outliers locate the tails of the distribution.
11. **Deep Support Vector Data Description (DeepSVDD)** [151]. 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)** [206]. 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 [77] optimizer is 0.0001 and is decayed once the number of epochs 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 are set to 0.1 and 0.005, respectively.
13. **Lightweight on-line detector of anomalies (LODA)** [136]. 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 online.
14. **Isolation Forest (IForest)** [100]. 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)** [7]. 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 with 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 [149] 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)** [152]. A deep one-class method that improves the unsupervised DeepSVDD [151] 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 the hypersphere, where the autoencoder 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 nEarrest Neighbor distance-based method (REPEN)** [127]. 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 [187] optimizer with 0.001 learning rate and 0.95 ρ is applied to update network parameters.

4. **Deviation Networks (DevNet)** [131]. A neural network-based model uses a prior probability to enforce a statistical 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 [149] 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)** [130]. 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 a learning rate of 0.001 and 0.01 weight decay is applied to update network parameters.
6. **Feature Encoding With Autoencoders for Weakly Supervised Anomaly Detection (FEAWAD)** [204]. A neural network-based model that incorporates the network architecture of DAGMM [206] with the deviation loss of DevNet [131]. 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)** [195]. XGBOD first uses the passed-in unsupervised outlier detectors to extract richer representations 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 to this augmented feature space. We use the default hyperparameters in PyOD.

(iii) 9 supervised algorithms:

1. **Naive Bayes (NB)** [14]. 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)** [31]. SVM is effective in high-dimensional spaces and could still be 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)** [148]. MLP uses the binary cross entropy loss to update network parameters.
4. **Random Forest (RF)** [21]. RF is a meta estimator that fits several 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)** [29]. 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)** [74]. 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)** [138]. 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)** [56]. This method introduces a ResNet-like architecture [62] for tabular based data. ResNet is trained for 100 epochs with 64 batch size. AdamW [108] optimizer with 0.001 learning rate is applied to update network parameters.
9. **Feature Tokenizer + Transformer (FTTransformer)** [56]. FTTransformer is an effective adaptation of the Transformer architecture [171] 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.

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/>

B.2 ADBench Dataset List

Overview. As described in §3.2, ADBench is the largest AD benchmark with 57 datasets. More specifically, Table B1 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 show 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 due to the computational cost. Fig. B1 provides the anomaly ratio distribution of the datasets, where the median is equal to 5%. We release all the datasets and their raw version(s) when possible at <https://github.com/Minqi824/ADBench/tree/main/datasets>

Newly-added Datasets in ADBench. Since most of the public datasets are relatively small and simple, we introduce 10 more complex datasets from CV and NLP domains with more samples and richer features in ADBench (highlighted in Table B1 in blue).

Reasoning of Using CV/NLP Datasets. It is often challenging to directly run large CV and NLP datasets on selected shallow methods, e.g., OCSVM [157] and kNN [144] with high time complexity, we follow DeepSAD [152], ADIB [33], and DATE [115] to extract representations of CV and NLP datasets by neural networks for downstream detection tasks. More specifically, ADIB [33] shows that “transferring features from semantic tasks can provide powerful and generic representations for various AD problems”, which is true even when the pre-trained task is only loosely related to downstream AD tasks. Similarly, DeepSAD [152] uses pre-trained autoencoder to extract features for training classical AD detectors like OCSVM [157] and IForest [100]. For NLP datasets, DATE [115] uses fastText [72] and Glove [134] embeddings for evaluating classical AD methods (e.g., OCSVM [157] and IForest [100]) against proposed methods in NLP datasets.

We want to elaborate further on the reasons for adapting CV and NLP datasets for tabular AD. First, some shallow models, such as OCSVM [157], cannot directly run on (large, high-dimensional) CV datasets. Second, it is interesting to see whether tabular AD methods can work on CV/NLP data representations, which carry values in real-world applications where deep models are infeasible to run. Moreover, the extracted representations often lead to better downstream detection results [33]. Thus, we extract features from CV and NLP datasets by deep models to create “tabular” versions of them. Although not perfect, this may provide insights into shallow methods’ performance on (originally infeasible) CV and NLP datasets.

CV Datasets: For MNIST-C, we set original MNIST images to be normal and corrupted images in MNIST-C to be abnormal, like in the recent work [9]. For MVTec-10, we test different AD algorithms on the 15 image sets, where anomalies correspond to various manufacturing defects. For CIFAR10, FashionMNIST, and SVHN, we follow previous works [151] [152] and set one of the multi-classes as normal and downsample the remaining classes to 5% of the total instances as anomalies by default, and report the average results over all the respective classes.

NLP Datasets: For Amazon and Imdb, we regard the original negative class as the anomaly class. For Yelp, we regard the reviews of 0 and 1 stars as the anomaly class, and the reviews of 3 and 4 stars as the normal class. For 20newsgroups dataset, like in DATE [115] and CVDD [153], we only take the articles from the six top-level classes: *computer, recreation, science, miscellaneous, politics, religion*. Similarly, for the multi-classes datasets 20newsgroups and Agnews, we set one of the classes as normal and downsample the remaining classes to 5% of the total instances as anomalies.

Backbone Choices of Feature Extraction. Pretrained models are applied to extract data embedding from CV and NLP datasets to access these more complex representations. For CV datasets, following [16] and [147], we use ResNet18¹ [62] pretrained on the ImageNet [35] to extract meaningful embedding after the last average pooling layer. We also provide the embedding version that are extracted by the ImageNet-pretrained ViT² [39]. For NLP datasets, instead of using traditional embedding methods like fastText [19] [72] or Glove [134], we apply BERT³ [75] pretrained on the BookCorpus and English Wikipedia to extract more enriching embedding of the [CLS] token. In addition, we provide the embedding version that are extracted by the pretrained RoBERTa⁴ [105]

¹https://pytorch.org/hub/pytorch_vision_resnet/

²<https://github.com/lukemelas/PyTorch-Pretrained-ViT>

³<https://huggingface.co/bert-base-uncased>

⁴<https://huggingface.co/roberta-base>

in our codebase^[1]. Although we release all the generated datasets for completeness, we analyze the results based on the datasets generated by BERT and ResNet18. Future work may consider analyzing the impact of backbones on detection performance.

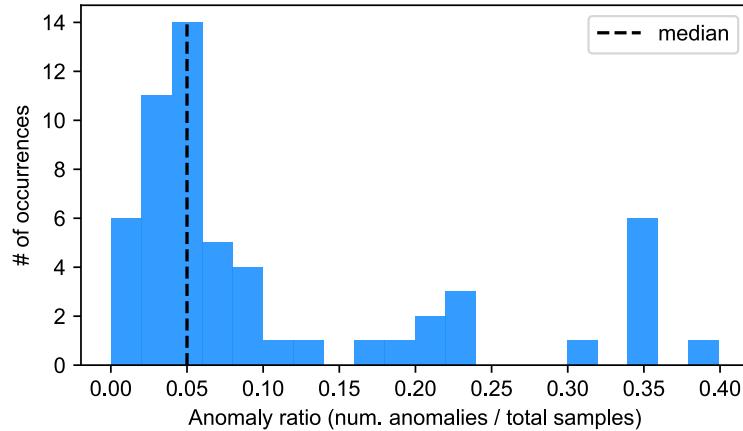


Figure B1: Distribution of anomaly ratios in 57 datasets in ADBench, where 40 datasets' anomaly ratio is below 10% (median=5%).

^[1]<https://github.com/Minqi824/ADBench/tree/main/datasets>

Table B1: Data description of the 57 datasets included in ADBench; 10 newly added datasets from CV and NLP domain are highlighted in blue at the bottom of the table.

Data	# Samples	# Features	# Anomaly	% Anomaly	Category	Reference
ALOI	49534	27	1508	3.04	Image	[42]
anthyroid	7200	6	534	7.42	Healthcare	[141]
backdoor	95329	196	2329	2.44	Network	[119]
breastw	683	9	239	34.99	Healthcare	[173]
campaign	41188	62	4640	11.27	Finance	[131]
cardio	1831	21	176	9.61	Healthcare	[12]
Cardiotocography	2114	21	466	22.04	Healthcare	[12]
celeba	202599	39	4547	2.24	Image	[131]
census	299285	500	18568	6.20	Sociology	[131]
cover	286048	10	2747	0.96	Botany	[18]
donors	619326	10	36710	5.93	Sociology	[131]
fault	1941	27	673	34.67	Physical	[42]
fraud	284807	29	492	0.17	Finance	[131]
glass	214	7	9	4.21	Forensic	[43]
Hepatitis	80	19	13	16.25	Healthcare	[36]
http	567498	3	2211	0.39	Web	[145]
InternetAds	1966	1555	368	18.72	Image	[25]
Ionosphere	351	33	126	35.90	Oryctognosy	[163]
landsat	6435	36	1333	20.71	Astronautics	[42]
letter	1600	32	100	6.25	Image	[48]
Lymphography	148	18	6	4.05	Healthcare	[26]
magic.gamma	19020	10	6688	35.16	Physical	[42]
mammography	11183	6	260	2.32	Healthcare	[176]
mnist	7603	100	700	9.21	Image	[90]
musk	3062	166	97	3.17	Chemistry	[37]
optdigits	5216	64	150	2.88	Image	[10]
PageBlocks	5393	10	510	9.46	Document	[113]
pendigits	6870	16	156	2.27	Image	[9]
Pima	768	8	268	34.90	Healthcare	[145]
satellite	6435	36	2036	31.64	Astronautics	[145]
satimage-2	5803	36	71	1.22	Astronautics	[145]
shuttle	49097	9	3511	7.15	Astronautics	[145]
skin	245057	3	50859	20.75	Image	[42]
smtp	95156	3	30	0.03	Web	[145]
SpamBase	4207	57	1679	39.91	Document	[25]
speech	3686	400	61	1.65	Linguistics	[23]
Stamps	340	9	31	9.12	Document	[25]
thyroid	3772	6	93	2.47	Healthcare	[142]
vertebral	240	6	30	12.50	Biology	[17]
vowels	1456	12	50	3.43	Linguistics	[82]
Waveform	3443	21	100	2.90	Physics	[107]
WBC	223	9	10	4.48	Healthcare	[114]
WDBC	367	30	10	2.72	Healthcare	[114]
Wilt	4819	5	257	5.33	Botany	[25]
wine	129	13	10	7.75	Chemistry	[2]
WPBC	198	33	47	23.74	Healthcare	[114]
yeast	1484	8	507	34.16	Biology	[66]
CIFAR10	5263	512	263	5.00	Image	[81]
FashionMNIST	6315	512	315	5.00	Image	[178]
MNIST-C	10000	512	500	5.00	Image	[120]
MVTec-AD	See Table B2.				Image	[16]
SVHN	5208	512	260	5.00	Image	[121]
Agnews	10000	768	500	5.00	NLP	[192]
Amazon	10000	768	500	5.00	NLP	[63]
Imdb	10000	768	500	5.00	NLP	[111]
Yelp	10000	768	500	5.00	NLP	[192]
20newsgroups	See Table B3.				NLP	[86]

Table B2: Detailed description of the MVTec-AD dataset; see the full dataset list in Table B1. For MVTec-AD dataset, we evaluate 30 algorithms on each class and report the average performance of all classes.

Class	# Samples	# Features	# Anomaly	% Anomaly
Carpet	397	512	89	22.42
Grid	342	512	57	16.67
Leather	369	512	92	24.93
Tile	347	512	84	24.21
Wood	326	512	60	18.40
Bottle	292	512	63	21.58
Cable	374	512	92	24.60
Capsule	351	512	109	31.05
Hazelnut	501	512	70	13.97
Metal Nut	335	512	93	27.76
Pill	434	512	141	32.49
Screw	480	512	119	24.79
Toothbrush	102	512	30	29.41
Transistor	313	512	40	12.78
Zipper	391	512	119	30.43
Total	5354	512	1258	23.50

Table B3: Detailed description of the 20newsgroups dataset; see the full dataset list in Table B1. For 20newsgroups dataset, we evaluate 30 algorithms on each class and report the average performance of all classes.

Class	# Samples	# Features	# Anomaly	% Anomaly
Computer	3090	768	154	4.98
Recreation	2514	768	125	4.97
Science	2497	768	124	4.97
Miscellaneous	615	768	30	4.88
Politics	1657	768	82	4.95
Religion	1532	768	76	4.96
Total	11905	768	591	4.96

B.3 Additional Demonstration of Synthetic Anomalies for §3.3.2

In addition to Fig. 3 that demonstrates the synthetic anomalies on Lymphography dataset in §3.3.2, we provide another example here for Ionosphere data.

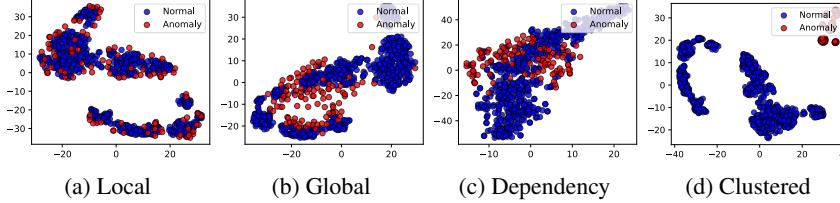


Figure B2: Illustration of four types of synthetic anomalies shown on Ionosphere dataset.

B.4 Open-source Release

As mentioned before, the full experiment code, datasets, and 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 a newly proposed AD algorithm easily.

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 newcomer data), we adapt all the algorithms for the inductive setting to predict the newcomer data, which is helpful in applications and also common in popular AD library PyOD [198], TODS [84] [85], and PyGOD [102]. Thus, we use 70% data for training and the remaining 30% as a testing set. We use stratified sampling to keep the anomaly ratio consistent. We repeat each experiment 3 times and report the average. The 10 complex CV and NLP datasets are mainly considered for evaluating algorithm performance on the public 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 noise 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). Future works may resort to the help of the latest generative methods like diffusion models [184].

Hyperparameter Settings. For all the algorithms in ADBench, we use their default hyperparameter (HP) settings in the original paper for a fair comparison. Specific values can be found in Appx B.1 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 [164], while we do not consider this setting in this work.

Extensive Experiments. In total ADBench conducts 98,436 experiments, where each denotes one algorithm’s result on a dataset under a specific setting. More specifically, we have 27,090 experiments in §4.2. For 47 classical datasets:

- Unsupervised methods on benchmark `real-world` datasets {14 algorithms, 47 datasets, 3 repeat times} leads to 1,974 experiments.
- Semi- and fully-supervised on real-world datasets {16 algorithms, 47 datasets, 3 repeat times, 7 settings of labeled anomalies} leads to 15,792 experiments.

As we described in Appx. B.2 we totally have 74 subclasses for the 10 CV and NLP datasets, thus generating:

- Unsupervised methods on benchmark datasets {14 algorithms, 74 subclasses, 1 repeat times} leads to 1,036 experiments.
- Semi- and fully-supervised on real-world datasets {16 algorithms, 74 subclasses, 1 repeat times, 7 settings of labeled anomalies} leads to 8,288 experiments.

Additionally, we have 17,766 experiments for understanding the algorithm performances under four types of anomalies in §4.3:

- Unsupervised methods on benchmark `real-world` datasets {14 algorithms, 47 datasets, 3 repeat times} leads to 1,974 experiments.
- Semi- and fully-supervised on benchmark datasets {16 algorithms, 47 datasets, 3 repeat times, 7 settings of labeled anomalies} leads to 15,792 experiments.

Finally, we have 53,580 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, 47 datasets, 3 repeat times, 5 settings of data noises, 2 scenarios} leads to 42,300 experiments.
- For annotation errors {16 algorithms, 47 datasets, 3 repeat times, 5 settings of data noises} leads to 11,280 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 experiments on an NVIDIA Tesla V100 GPU accelerator. The model runtime on benchmark datasets is reported in Appx. D.1

¹ ADBench repo: <https://github.com/Minqi1824/ADBench>

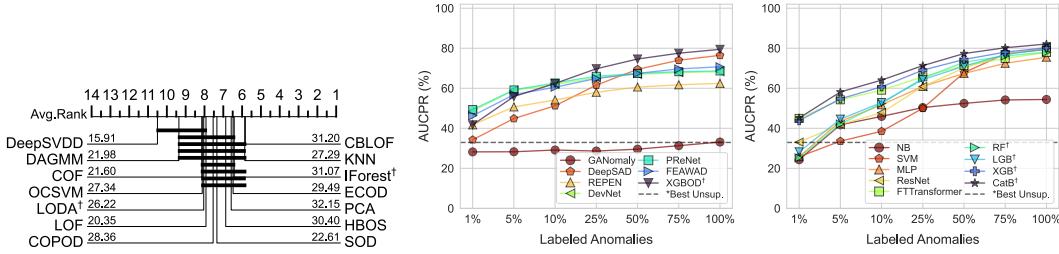
D Additional Experiment Results

D.1 Additional Results for Overall Model Performance on Benchmark Datasets in §4.2

In addition to the AUCROC results presented in §4.2 we also show the AUCPR results of model performance on 57 benchmark datasets in Fig. D3, 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 models leveraging labeled anomalies (see Fig. D3(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 [128] [130] [131].

Fig. D4 and D5 show the boxplots of AUCROC and AUCPR of 30 algorithms on the 57 benchmark datasets. These results validate the no-free-lunch theorem, where no model 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.

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



(a) Avg. rank (lower the better) and avg. AUCPR (on each line) of unsupervised methods; groups of algorithms not statistically different are connected horizontally. (b) Avg. AUCPR (on 57 datasets) vs. % of labeled anomalies (x-axis); semi-supervised (left) and fully-supervised (right). Most label-informed algorithms outperform the best *unsupervised* algorithm CBLOF (denoted as the dashed line) with 10% labeled anomalies.

Figure D3: AD model's AUCPR on 57 benchmark 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 ratios of labeled anomalies γ_l . The semi-supervised methods leverage the labels more efficiently w/ small γ_l .

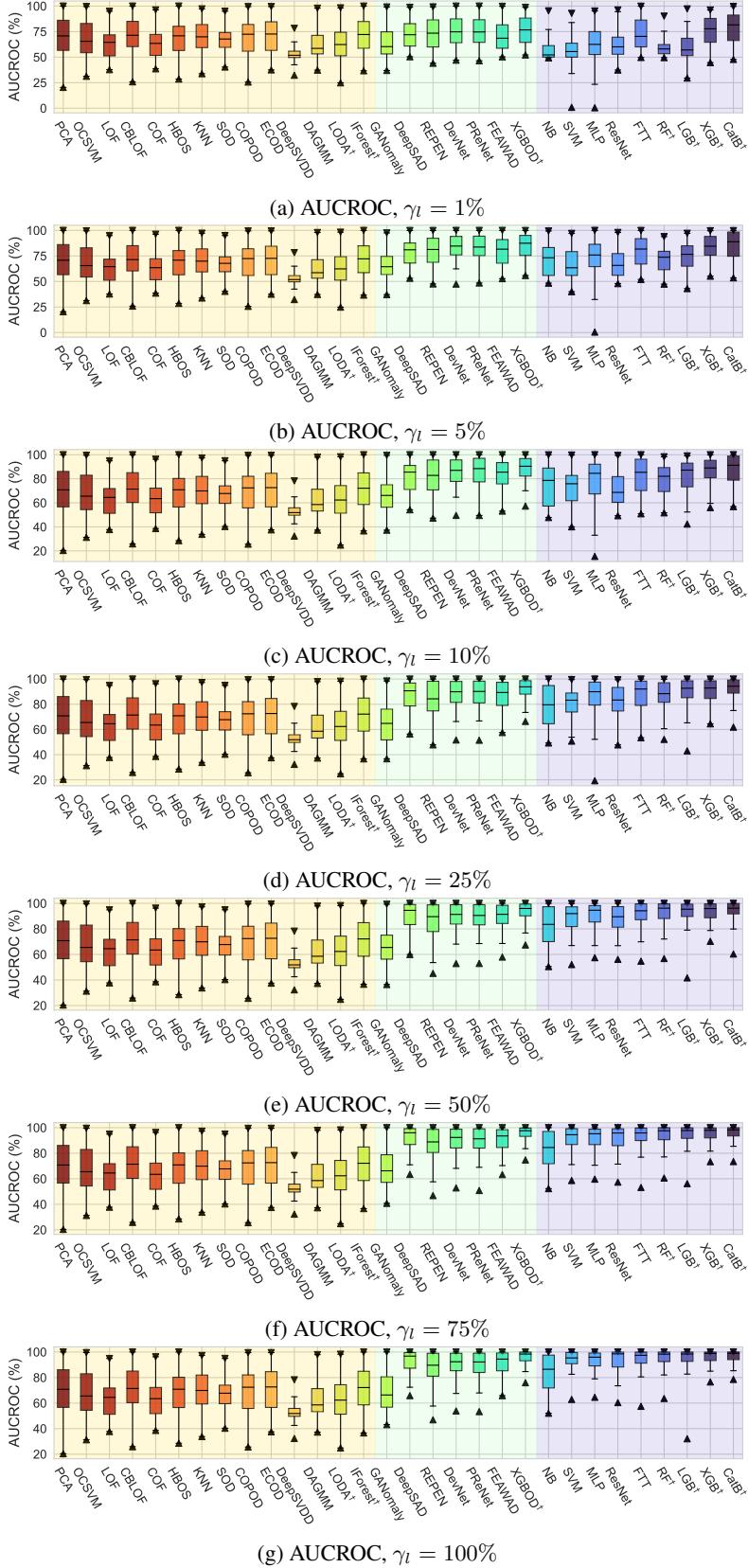


Figure D4: Boxplot of AUCROC. We denote unsupervised methods in light yellow, semi-supervised methods in light green, and supervised methods in light purple. Consistent with the CD diagrams, we notice that none of the unsupervised methods visually outperform.

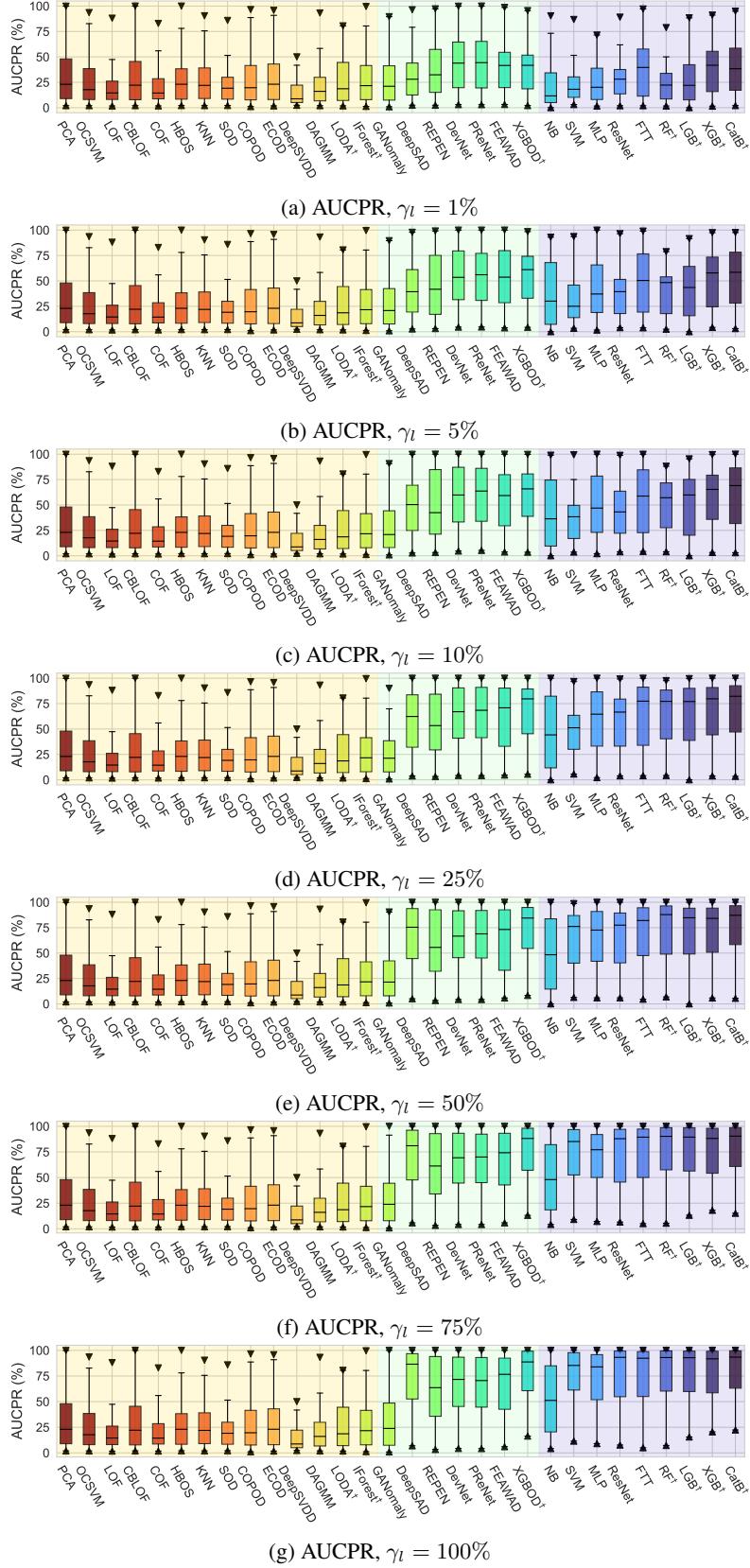


Figure D5: Boxplot of AUCPR. We denote unsupervised methods in light yellow, semi-supervised methods in light green, and supervised methods in light purple. Consistent with the CD diagrams, we notice that none of the unsupervised methods visually outperform.

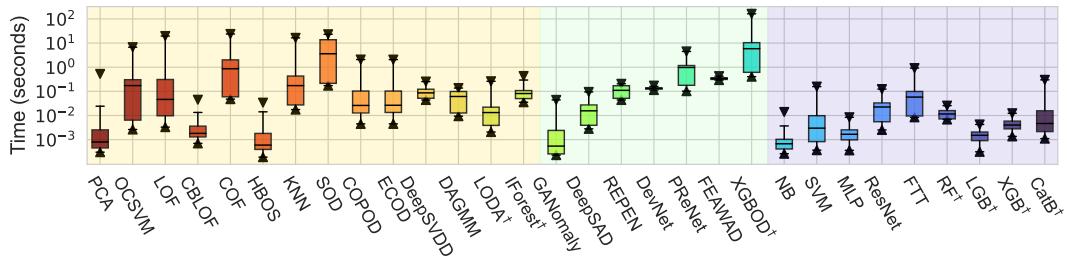


Figure D6: Inference time of included algorithms. We denote unsupervised methods in yellow (light yellow), semi-supervised methods in green (light green), and supervised methods in purple (light purple). Consistent with the train time in Fig. 4d PCA, HBOS, GANomaly and NB take the least inference time on test datasets, while more complex feature representation methods like SOD and XGBOD spend more time due to the search of the feature subspace.

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. D7 and Fig. D8, 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 to or even outperforms the best unsupervised model when more labeled anomalies are available. Recall that XGBOD employs the stacking ensemble method [174], 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 future research directions.

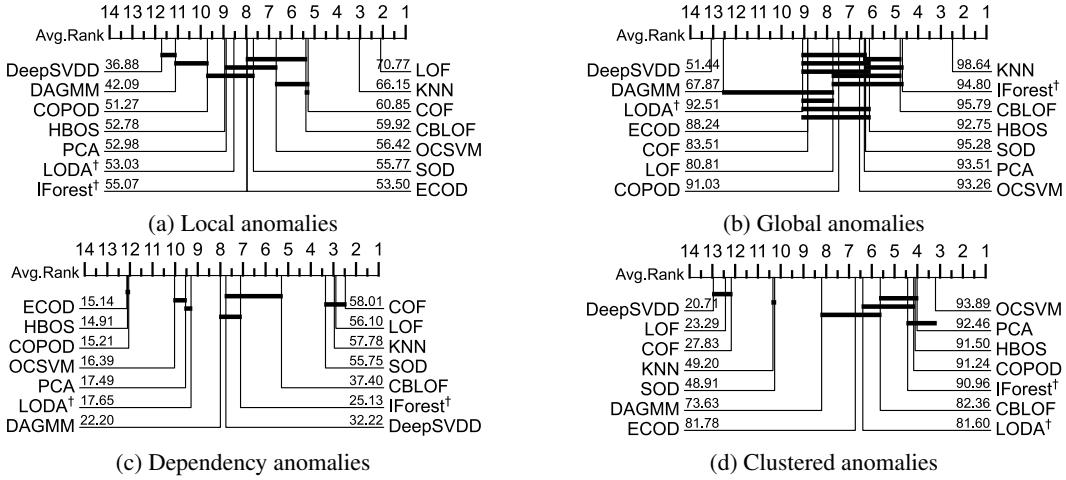


Figure D7: AUCPR CD Diagram of unsupervised methods on different types of anomalies. The unsupervised methods perform well when their assumptions conform to the anomaly types.

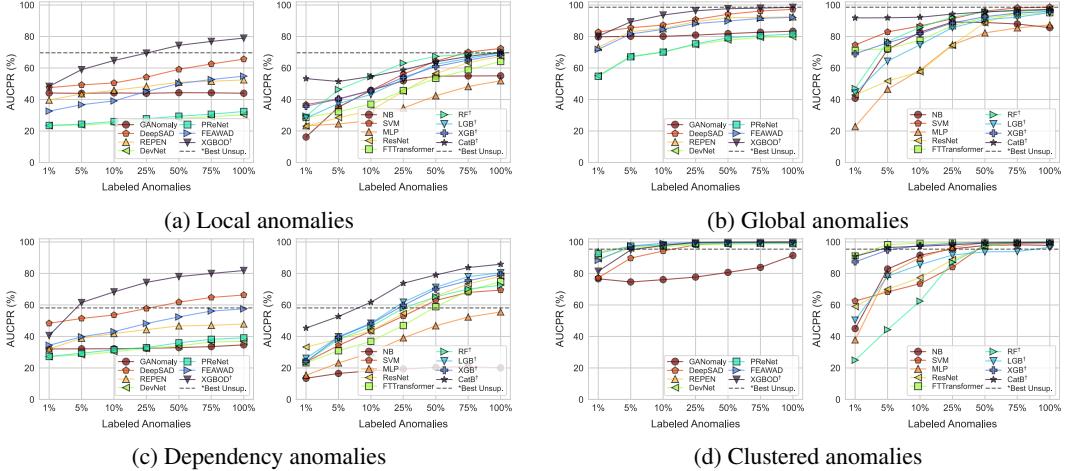


Figure D8: Semi- (left of each subfigure) and supervised (right) algorithms' performance on different types of anomalies with varying levels of labeled anomalies for AUCPR performance. Surprisingly, these label-informed algorithms are *inferior* to the best unsupervised method except for the clustered anomalies.

D.3 Additional Results for Algorithm Robustness in §4.4

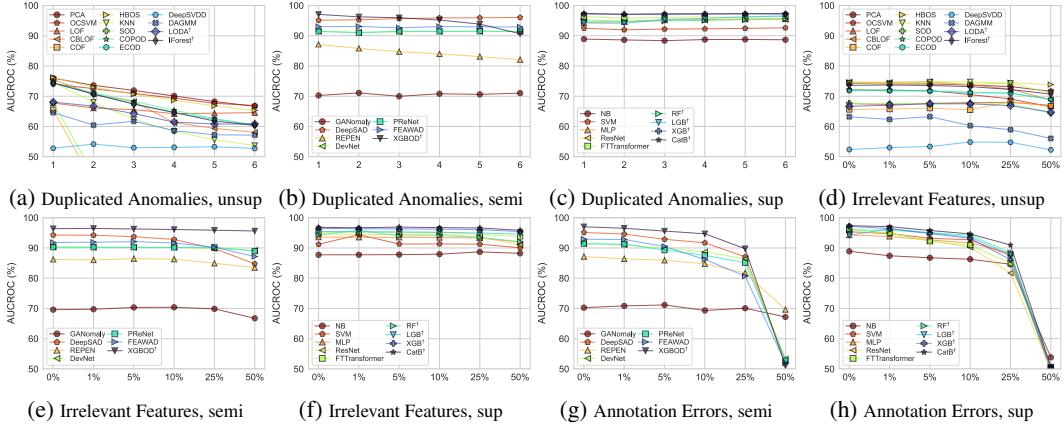


Figure D9: Algorithm performance under noisy and corrupted data (i.e., duplicated anomalies for (a)-(c), irrelevant features for (d)-(f), and annotation errors for (g) and (h)). X-axis denotes either the duplicated times or the noise ratio. Y-axis denotes the AUCROC performance and its range remains consistent across different algorithms. The results reveal 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. The results are mostly consistent with the observations in Fig. 7 (§4.4) showing the relative performance change.

In Fig. D9 we provide the performance of the AD algorithms under noisy and corrupted data. Along with the relative performance changes shown in Fig. 7 the analysis in §4.4 still stands.

In addition to the primary results shown in §4.4 we provide the AUCPR results for algorithm robustness in Fig. D10 and D11. 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. D10(a)-(c)). This is expected as AUCPR emphasizes the positive classes (i.e., anomalies), and more duplicated anomalies favor this metric. Since this observation is consistently true for both unsupervised and label-informed methods, it would not largely impact our selection of algorithms. However, if we care about both anomaly and normal classes equally, the results on AUCROC in §4.4 still stand —unsupervised methods are more susceptible to duplicate anomalies.

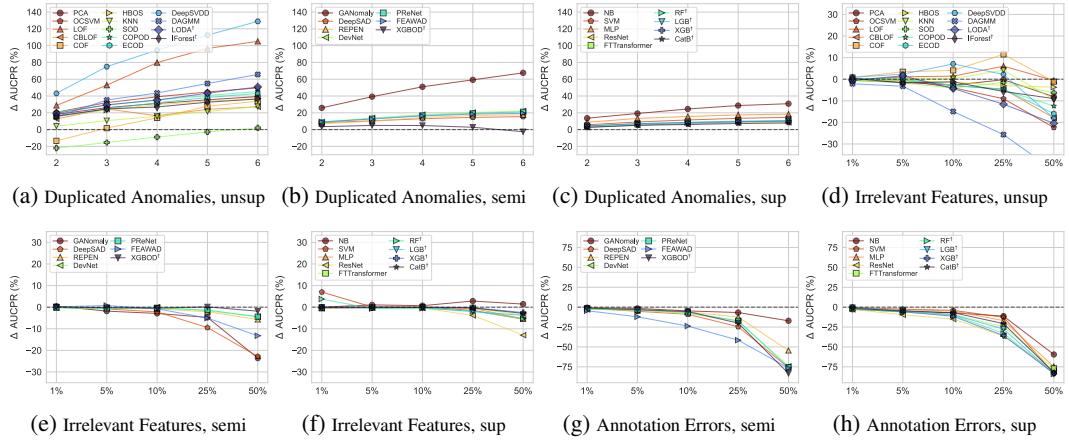


Figure D10: Algorithm performance change under noisy and corrupted data (i.e., duplicated anomalies for (a)-(c), irrelevant features for (d)-(f), and annotation errors for (g) and (h)). y-axis denotes the % of performance change (ΔAUCPR) and its range remains consistent across different algorithms. The 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) showing the AUCROC.

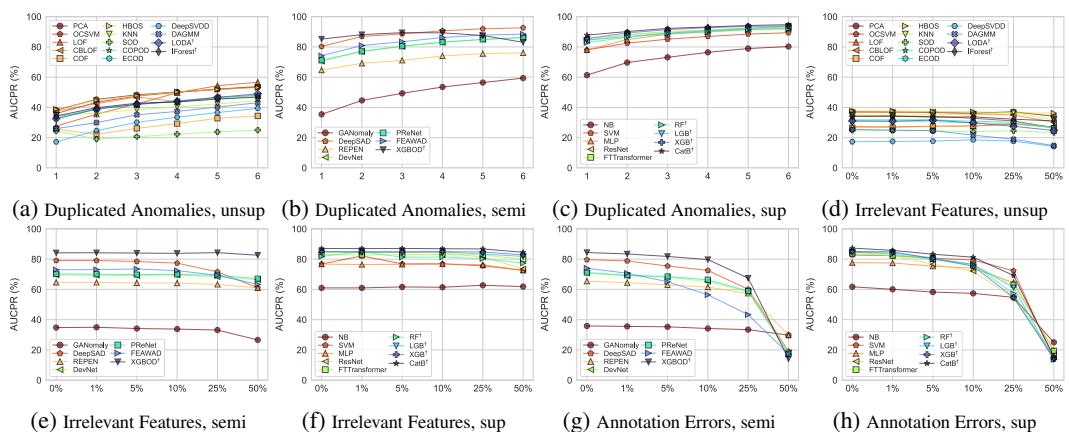


Figure D11: Algorithm performance under noisy and corrupted data (i.e., duplicated anomalies for (a)-(c), irrelevant features for (d)-(f), and annotation errors for (g) and (h)). X-axis denotes either the duplicated times or the noise ratio. Y-axis denotes the AUCPR performance and its range remains consistent across different algorithms. The results reveal 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.

D.4 Full Performance Tables on Benchmark 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 [199], where the historical algorithm performance table serves as a great source for building strong meta-learning methods.

Table D4: AUCROC of 14 unsupervised algorithms on 57 benchmark datasets. We show the performance rank in parenthesis (the lower, the better), and mark the best performing method(s) in **bold**.

Datasets	PCA	OCSVM	LOF	CBLOF	COF	HBOS	KNN	SOD	COPOD	ECOD	Deep SVDD	DA GMM	LODA	IForest	
ALOI	56.65(6)	55.85(8)	66.63(1)	55.22(9)	64.68(2)	52.63(11)	61.47(3)	61.09(4)	53.75(10)	56.60(7)	50.29(14)	51.96(12)	51.33(13)	56.66(5)	
amnthyroid	66.25(8)	57.23(12)	70.20(7)	62.26(10)	65.92(9)	60.15(11)	71.69(6)	77.38(3)	76.80(4)	78.03(2)	76.62(5)	56.53(13)	41.02(14)	82.01(1)	
backdoor	80.13(7)	86.20(2)	85.68(3)	81.16(4)	73.03(8)	71.43(10)	80.82(6)	69.54(11)	80.97(5)	86.33(1)	55.16(14)	56.26(13)	69.22(12)	72.15(9)	
breastw	95.13(8)	80.30(10)	46.61(12)	38.84(3)	98.94(3)	97.01(6)	99.68(1)	99.17(2)	65.66(1)	N/A(N/A)	98.49(4)	98.32(5)			
campaign	72.78(4)	65.52(9)	58.85(10)	66.61(8)	57.26(11)	76.81(1)	72.10(5)	69.04(7)	77.69(2)	76.78(3)	48.70(4)	56.08(12)	51.43(13)	71.71(6)	
cardio	95.55(1)	93.91(3)	66.33(13)	89.93(7)	71.41(9)	84.67(8)	76.64(9)	73.25(1)	92.35(5)	94.44(2)	58.96(4)	75.01(10)	90.34(6)	93.19(4)	
Cardiotocography	74.67(2)	77.86(1)	59.51(10)	64.54(7)	57.77(12)	60.86(9)	56.23(11)	51.69(14)	67.02(6)	68.92(4)	53.53(3)	62.01(8)	73.65(3)	67.57(5)	
celeba	79.38(1)	70.70(6)	38.55(14)	73.99(4)	38.58(13)	76.18(2)	59.63(9)	47.85(1)	75.68(3)	72.82(5)	50.36(10)	44.74(12)	60.11(8)	70.41(7)	
census	68.74(2)	54.58(10)	47.19(12)	59.41(8)	41.35(13)	64.94(5)	66.75(4)	62.31(6)	69.07(1)	68.44(3)	51.07(11)	59.29(9)	36.86(14)	59.52(7)	
cover	93.73(1)	92.62(3)	84.58(10)	89.30(6)	76.91(12)	80.24(11)	85.97(9)	74.46(13)	88.64(7)	93.42(2)	46.20(4)	89.89(5)	92.34(4)	86.74(8)	
donors	83.15(1)	71.93(7)	55.49(11)	60.44(10)	70.54(9)	78.23(4)	81.09(3)	55.21(12)	81.76(2)	74.45(6)	50.27(13)	70.57(8)	24.86(4)	77.68(5)	
fault	46.02(10)	47.69(9)	58.93(5)	64.06(3)	62.10(4)	51.28(8)	72.98(1)	68.11(2)	43.88(12)	43.41(13)	51.67(7)	45.86(11)	41.71(14)	57.02(6)	
fraud	90.35(8)	90.62(6)	94.92(2)	91.70(5)	93.05(4)	90.29(9)	93.56(3)	94.97(1)	88.32(13)	89.85(10)	64.98(4)	89.53(11)	88.99(12)	90.38(7)	
glass	66.29(12)	35.36(14)	69.20(11)	82.94(1)	72.24(10)	77.23(3)	82.29(2)	73.36(7)	72.43(9)	75.70(6)	47.49(3)	76.09(5)	73.13(8)	77.13(4)	
Hepatitis	75.95(4)	67.75(7)	38.02(14)	66.40(8)	41.45(13)	79.85(2)	52.76(11)	68.17(6)	82.05(1)	79.67(3)	50.96(2)	54.80(10)	64.87(9)	69.75(5)	
http	99.72(2)	99.59(4)	27.46(11)	99.60(3)	88.78(8)	99.53(5)	3.37(13)	78.04(9)	99.29(6)	98.10(7)	69.05(10)	N/A(N/A)	12.48(2)	99.96(1)	
InternetAds	61.67(11)	68.28(4)	65.83(8)	70.58(1)	63.79(9)	68.03(5)	69.99(2)	61.85(10)	67.05(7)	67.10(6)	60.20(12)	N/A(N/A)	55.38(13)	69.01(3)	
Ionosphere	79.19(8)	75.92(10)	90.32(2)	90.72(1)	86.76(4)	62.49(13)	88.26(3)	86.23(5)	79.34(7)	75.59(1)	50.89(4)	73.41(12)	78.42(9)	84.50(6)	
landsat	35.76(14)	36.15(13)	53.90(7)	63.55(2)	53.50(8)	55.14(6)	57.95(4)	59.54(3)	41.55(11)	56.61(5)	63.61(1)	43.92(10)	38.17(12)	47.64(9)	
letter	50.29(12)	46.18(14)	84.49(2)	75.62(5)	80.03(4)	59.74(7)	56.69(4)	86.19(1)	84.09(3)	54.32(9)	50.76(10)	56.64(8)	50.42(11)	50.24(13)	61.07(6)
Lymphography	99.82(2)	99.47(4)	89.86(1)	99.83(1)	90.85(9)	99.49(6)	55.91(13)	72.49(1)	99.48(7)	99.52(5)	32.29(4)	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)	75.25(2)	68.33(7)	64.36(11)	60.46(13)	58.58(14)	68.22(8)	73.25(4)	
mammography	88.72(3)	84.95(6)	74.73(12)	83.77(9)	75.73(11)	86.27(5)	84.53(7)	81.51(11)	90.69(2)	90.75(1)	56.98(13)	N/A(N/A)	83.91(8)	86.39(4)	
mnist	85.29(1)	82.95(3)	67.13(11)	79.45(6)	70.78(9)	60.42(12)	80.58(5)	60.10(13)	77.74(7)	84.60(2)	53.40(4)	67.23(10)	72.27(3)	80.98(4)	
musik	100.00(1)	80.58(8)	41.18(13)	100.00(1)	38.66(4)	100.00(1)	69.01(11)	74.09(10)	94.42(7)	91.11(6)	43.52(12)	76.85(9)	95.11(5)	99.99(4)	
optdigits	51.65(11)	54.00(6)	56.10(9)	81.51(1)	49.15(12)	81.63(2)	41.73(13)	58.92(8)	66.71(4)	61.04(7)	38.89(11)	62.57(6)	61.74(6)	70.92(3)	
PageBlocks	99.64(2)	88.85(5)	70.90(12)	85.04(7)	72.65(13)	80.58(10)	81.94(9)	77.75(11)	88.05(6)	90.93(1)	57.77(4)	89.61(3)	83.34(8)	89.57(4)	
pendigits	93.73(3)	83.78(2)	72.99(12)	90.40(7)	45.07(13)	93.94(4)	72.95(9)	66.20(10)	90.68(6)	91.54(1)	39.92(4)	64.32(11)	89.10(8)	94.41(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)	69.10(6)	51.54(13)	51.03(4)	55.92(12)	69.93(8)	72.87(2)	
satellite	59.62(10)	59.02(11)	55.88(12)	71.32(3)	54.74(14)	74.80(2)	65.18(5)	63.96(6)	63.20(7)	75.06(1)	55.30(3)	62.33(8)	61.98(9)	70.43(4)	
satimage-2	97.62(4)	97.35(6)	47.36(14)	99.84(1)	56.70(12)	97.65(3)	92.60(10)	83.08(11)	97.21(7)	97.11(8)	53.14(3)	96.29(9)	97.56(5)	99.16(2)	
shuttle	98.63(5)	97.40(7)	57.11(12)	83.48(8)	51.77(4)	98.63(4)	69.64(9)	69.51(10)	99.35(3)	99.40(2)	52.05(3)	97.92(6)	60.95(1)	99.56(1)	
smrip	45.26(10)	49.45(6)	46.47(8)	69.49(2)	41.66(12)	60.15(5)	71.46(1)	60.35(4)	47.55(7)	39.09(13)	44.05(1)	N/A(N/A)	45.75(9)	68.21(3)	
spambase	88.41(3)	70.80(4)	71.84(10)	79.68(5)	70.60(6)	70.52(12)	59.86(14)	79.09(7)	71.86(9)	78.24(8)	71.32(11)	67.43(3)	89.73(1)		
speech	54.66(6)	52.47(9)	43.33(11)	54.97(5)	40.96(13)	64.74(4)	53.38(5)	62.35(10)	70.09(1)	66.89(2)	53.55(8)	N/A(N/A)	41.99(12)	64.76(3)	
Stamps	50.79(9)	50.19(13)	52.48(6)	50.58(2)	55.97(1)	50.50(1)	51.03(8)	55.86(2)	52.89(4)	51.88(7)	53.43(3)	52.75(6)	49.84(14)	50.74(16)	
thyroid	91.47(2)	83.86(8)	51.26(14)	68.18(11)	53.81(13)	90.73(5)	68.61(10)	73.26(9)	93.40(1)	91.41(3)	55.84(2)	88.88(6)	87.18(7)	91.21(4)	
vertebral	96.34(3)	87.92(10)	86.86(11)	94.73(6)	90.87(9)	95.62(5)	95.93(4)	92.81(8)	94.30(7)	97.78(2)	49.64(1)	79.75(12)	74.30(13)	98.30(1)	
vorwells	37.06(8)	37.99(6)	49.29(2)	41.41(4)	48.71(3)	28.56(13)	33.79(11)	40.32(5)	25.64(14)	37.51(7)	36.67(9)	53.20(1)	30.57(12)	36.66(10)	
WBC	98.20(7)	99.03(4)	54.17(14)	99.46(1)	71.44(1)	78.68(10)	86.60(2)	81.73(7)	81.07(8)	86.57(1)	54.47(13)	49.35(14)	60.13(11)	71.47(7)	
WDBC	99.05(4)	98.86(6)	89.00(12)	99.32(3)	96.26(9)	99.50(1)	91.72(11)	91.90(10)	99.42(2)	97.20(8)	65.69(14)	76.67(13)	98.26(7)	98.95(5)	
Wilt	20.39(14)	31.28(12)	50.65(2)	32.54(10)	49.66(3)	32.49(11)	48.42(4)	53.25(1)	33.40(9)	39.43(7)	46.08(5)	37.29(8)	26.42(13)	41.94(6)	
wine	84.37(4)	37.07(6)	37.74(13)	25.86(4)	44.44(12)	91.36(1)	44.98(11)	46.11(10)	88.65(3)	71.34(7)	59.52(9)	61.70(8)	90.12(2)	80.37(5)	
WPBC	46.01(10)	45.35(12)	41.41(14)	44.77(13)	48.88(11)	51.24(1)	46.59(9)	51.14(2)	49.34(4)	46.83(7)	49.79(3)	47.80(6)	49.31(5)	46.63(8)	
yeast	41.15(7)	41.00(9)	45.31(2)	44.85(3)	44.48(5)	39.64(10)	39.06(12)	42.46(6)	36.99(14)	39.61(11)	47.92(1)	41.11(8)	44.58(4)	37.76(13)	
CIFAR10	63.87(6)	63.76(7)	68.57(1)	64.23(4)	64.70(3)	57.50(13)	64.75(2)	64.22(5)	58.64(11)	61.04(10)	56.04(12)	58.08(12)	62.34(8)	61.28(9)	
FashionMNIST	86.09(3)	85.24(4)	67.57(12)	88.17(1)	71.44(1)	78.68(10)	86.60(2)	81.73(7)	81.07(8)	83.63(6)	63.32(14)	67.29(13)	80.28(9)	84.89(5)	
MNIST-C	73.75(5)	72.21(8)	68.27(12)	80.86(2)	69.81(11)	70.82(10)	81.26(1)	74.00(4)	71.26(9)	72.64(7)	51.85(14)	58.56(13)	74.37(3)	73.74(6)	
MVTec-AD	72.42(8)	69.84(10)	74.19(2)	75.98(1)	69.70(11)	73.36(4)	72.96(6)	71.57(9)	72.91(7)	73.46(3)	57.10(14)	66.47(13)	68.51(12)	73.19(5)	
SVHN	60.53(6)	60.73(5)	64.51(1)	60.30(7)	63.47(2)	56.08(13)	62.63(3)	61.09(4)	56.75(12)	58.27(9)	53.47(14)	57.22(11)	58.26(10)	58.62(8)	
Agnews	54.70(8)	54.34(9)	71.80(1)	60.02(5)	68.97(2)	53.87(10)	64.11(3)	62.81(4)	52.98(12)	53.04(11)	42.51(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)	60.65(1)	56.94(5)	56.79(6)	39.08(14)	53.58(13)	54.20(11)	56.13(8)	
Indb	47.06(12)	46.07(14)	48.71(9)	49.35(6)	49.64(5)	49.10(7)	47.83(10)	49.86(4)	50.68(3)	50.73(1)	50.73(2)	47.67(11)	46.43(13)	49.09(8)	
Yelp	60.71(11)	60.28(12)	67.09(3)	64.90(5)	66.11(4)	61.85(9)	69.84(1)	67.74(2)	62.36(7)	62.15(8)	54.62(14)	56.28(13)	61.36(10)	62.53(6)	
2news	56.66(7)	56.45(8)	62.14(1)	57.59(5)	61.80(2)	56.28(9)	59.33(3)	58.56(4)	55.79(11)	56.00(10)	50.24(14)	54.17(13)	55.53(12)	56.90(6)	

Table D5: AUCPR of 14 unsupervised algorithms on 57 benchmark datasets. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	PCA	OCSVM	LOF	CBLOF	COF	HBOS	KNN	SOD	COPOD	ECOD	Deep SVDD	DA GMM	LODA	IForest
ALOI	4.17(9)	5.02(5)	8.08(1)	4.46(7)	6.85(2)	3.69(13)	6.02(3)	5.97(4)	3.62(14)	3.90(11)	4.01(10)	4.33(8)	4.53(6)	3.90(12)
amnthyroid	16.12(8)	10.37(12)	15.71(9)	13.69(11)	14.39(10)	16.99(5)	16.74(6)	18.84(4)	16.58(7)	24.65(2)	21.95(3)	9.64(13)	7.06(14)	30.47(1)
backdoor	31.29(3)	9.69(9)	26.14(4)	6.96(11)	24.68(5)	4.91(13)	45.22(1)	39.41(2)	7.69(10)	11.25(8)	12.85(7)	6.50(12)	14.51(6)	4.75(14)
breastw	95.11(6)	82.70(10)	28.55(12)	91.54(8)	27.60(13)	97.71(3)	92.19(7)	84.88(9)	99.40(1)	98.54(2)	50.92(11)	N/A(N/A)	97.04(4)	96.04(5)
campagna	27.90(6)	29.22(5)	14.51(11)	23.99(8)	13.01(13)	37.99(2)	27.18(7)	18.88(9)	38.58(1)	37.40(3)	11.60(4)	14.62(10)	13.47(12)	32.26(4)
cardio	66.06(2)	62.89(3)	23.79(13)	61.95(4)	28.67(11)	52.10(8)	40.72(9)	28.54(2)	60.42(5)	68.59(1)	22.50(4)	28.92(10)	53.41(7)	59.95(6)
Cardiotocography	47.95(3)	52.61(1)	30.66(11)	45.44(4)	28.21(13)	38.28(8)	34.79(9)	27.99(14)	40.46(7)	43.57(5)	34.03(10)	30.61(12)	48.00(2)	41.47(6)
celeba	15.89(1)	10.73(6)	1.71(14)	11.33(5)	1.77(13)	13.82(2)	3.14(9)	2.66(10)	13.69(3)	12.37(4)	2.34(11)	1.95(12)	4.04(8)	8.96(7)
census	10.02(1)	6.76(11)	5.45(12)	7.44(9)	4.88(14)	8.69(6)	9.00(4)	8.52(7)	9.92(2)	9.72(3)	6.87(10)	8.71(5)	5.01(13)	7.78(8)
cover	9.80(6)	11.41(4)	8.12(8)	5.83(12)	4.00(13)	6.83(10)	6.16(11)	3.88(4)	11.37(5)	15.63(2)	8.12(9)	27.59(1)	13.06(3)	8.85(7)
donors	17.90(3)	9.86(8)	7.88(11)	6.89(12)	8.80(10)	23.36(1)	14.75(4)	9.69(9)	21.58(2)	14.17(5)	6.38(13)	10.53(7)	3.78(14)	12.74(6)
fault	32.76(11)	38.44(7)	38.38(8)	43.98(3)	41.56(4)	36.47(9)	54.45(1)	48.42(2)	30.54(14)	30.82(13)	39.15(6)	33.48(10)	31.03(2)	41.09(5)
fraud	22.91(10)	47.58(1)	47.40(3)	47.52(2)	22.86(11)	25.89(9)	47.30(8)	31.37(8)	42.82(7)	42.99(6)	8.97(14)	21.52(3)	46.37(5)	21.67(2)
glass	10.05(11)	8.04(8)	20.11(3)	13.84(6)	1.81(19)	11.82(8)	20.26(2)	18.73(4)	9.78(12)	18.43(5)	8.72(13)	24.50(1)	15.77(7)	10.99(10)
Hepatitis	36.65(4)	29.44(7)	13.69(14)	31.54(5)	14.39(13)	37.73(3)	21.95(2)	24.89(9)	41.50(1)	37.82(2)	22.77(11)	22.96(10)	30.90(6)	26.25(8)
http	56.43(2)	46.86(4)	3.82(11)	47.53(3)	9.57(9)	44.79(5)	0.70(12)	8.33(10)	35.19(6)	16.61(8)	N/A(N/A)	0.67(13)	90.83(1)	
InternetAds	32.55(10)	54.68(2)	40.49(8)	58.13(1)	38.67(9)	53.73(3)	43.23(7)	27.69(12)	50.97(5)	51.07(4)	27.91(11)	N/A(N/A)	23.89(3)	48.60(9)
Ionosphere	73.82(8)	74.54(7)	88.46(2)	90.22(2)	83.85(5)	41.76(7)	90.41(1)	85.87(4)	69.01(10)	65.60(11)	41.73(3)	64.98(12)	73.04(9)	80.16(6)
landsat	16.18(14)	16.21(13)	24.69(6)	30.97(2)	24.95(5)	22.03(9)	24.65(7)	26.38(3)	17.48(2)	25.17(4)	38.83(1)	24.48(8)	18.86(11)	19.81(10)
letter	6.46(12)	6.10(14)	14.80(5)	21.43(4)	8.38(9)	30.02(2)	26.63(3)	6.77(3)	9.94(10)	9.29(7)	11.68(6)	6.37(11)	8.49(8)	
Lymphography	97.03(3)	93.59(4)	23.38(11)	97.62(1)	36.68(9)	91.83(5)	38.69(9)	22.65(12)	88.68(7)	4.50(14)	19.52(3)	44.54(8)	97.31(2)	
magic_gamma	59.27(6)	51.43(12)	54.76(9)	68.85(2)	54.12(11)	75.63(1)	67.89(3)	59.79(7)	45.38(10)	49.17(13)	64.92(4)	44.49(8)	64.72(4)	
mnist	19.25(5)	12.94(6)	9.80(12)	11.14(10)	11.14(11)	2.13(3)	15.17(6)	13.41(8)	40.67(2)	4.26(13)	N/A(N/A)	14.75(7)	20.49(4)	
mnist	39.92(1)	33.20(3)	20.90(11)	28.82(5)	25.51(8)	12.51(14)	35.53(2)	19.15(13)	21.35(10)	41.93(4)	19.72(12)	23.75(9)	25.86(7)	27.71(6)
mnist	99.89(3)	10.61(9)	2.82(13)	100.00(2)	2.61(14)	100.00(1)	9.65(10)	7.59(9)	34.79(7)	34.95(6)	5.39(12)	32.78(8)	47.60(5)	99.61(4)
optdigits	2.76(13)	2.92(12)	6.06(3)	10.08(1)	4.42(6)	10.03(2)	3.06(11)	4.39(7)	4.36(8)	3.43(10)	2.50(14)	5.59(4)	3.05(9)	5.09(5)
PageBlocks	51.71(2)	49.14(6)	39.64(10)	49.65(4)	41.02(9)	33.32(3)	45.39(8)	37.83(11)	37.65(12)	49.30(5)	31.45(4)	53.25(1)	51.29(3)	46.04(7)
pendigits	23.65(3)	23.52(4)	3.78(12)	17.27(8)	2.89(13)	29.27(1)	6.50(9)	4.46(11)	21.22(6)	23.07(5)	2.45(14)	4.67(10)	15.71(7)	26.05(2)
Pima	54.03(5)	50.00(7)	47.18(9)	52.19(6)	47.70(10)	56.61(1)	55.14(4)	48.24(8)	55.19(3)	37.30(13)	35.87(7)	41.55(12)	44.09(11)	55.82(3)
satellite	59.64(6)	57.61(8)	37.68(14)	61.48(5)	39.70(13)	67.25(1)	50.01(10)	47.23(11)	56.58(9)	65.94(2)	40.11(12)	58.33(7)	61.94(4)	65.92(3)
satimage-2	85.69(3)	82.71(4)	4.29(13)	97.09(1)	8.80(12)	78.04(6)	39.14(9)	26.11(10)	76.55(7)	63.25(8)	3.08(14)	22.07(11)	80.52(5)	93.45(2)
shuttle	92.35(6)	85.29(7)	13.76(13)	60.98(8)	12.17(14)	96.40(3)	20.38(10)	20.27(11)	96.56(2)	95.76(4)	15.86(2)	93.20(5)	48.75(9)	97.62(1)
skin	17.40(11)	19.03(6)	18.25(9)	29.82(1)	16.38(12)	23.70(5)	28.72(2)	24.61(4)	17.99(10)	15.96(13)	18.48(7)	N/A(N/A)	18.44(8)	26.08(3)
smtip	66.70(2)	18.90(12)	20.69(11)	61.13(3)	35.20(8)	66.70(1)	33.36(9)	1.08(14)	51.01(6)	50.02(5)	50.03(4)	35.77(7)	1.24(13)	
SpamBase	41.57(6)	40.12(9)	35.16(12)	41.18(8)	34.73(13)	50.03(4)	41.42(7)	40.03(10)	56.68(1)	53.95(2)	42.23(3)	N/A(N/A)	35.88(11)	51.75(3)
speech	1.97(10)	1.96(11)	2.13(12)	1.99(9)	2.25(4)	2.09(6)	2.02(8)	2.13(5)	1.94(12)	1.77(14)	5.12(1)	2.03(7)	1.79(13)	2.31(3)
Stamps	41.09(3)	31.39(8)	21.29(11)	23.66(9)	16.50(13)	35.24(6)	23.52(10)	20.28(12)	43.10(2)	38.17(5)	11.40(14)	43.72(1)	34.69(7)	39.49(4)
thyroid	44.34(4)	21.23(9)	20.81(10)	29.95(6)	28.50(7)	50.98(5)	34.98(5)	23.56(8)	19.64(11)	54.05(2)	2.50(14)	16.06(12)	14.68(13)	63.11(1)
vertebral	10.49(10)	10.94(7)	14.24(2)	11.58(5)	13.85(3)	9.23(13)	10.57(8)	11.79(4)	8.89(4)	11.24(6)	10.49(9)	15.24(1)	9.68(12)	10.46(11)
vowels	8.92(10)	8.24(11)	34.42(4)	22.12(5)	55.96(2)	13.41(8)	63.41(1)	38.88(3)	4.14(13)	3.92(14)	4.99(12)	12.22(9)	13.82(7)	15.12(6)
Waveform	5.79(10)	4.37(13)	11.33(4)	18.98(1)	14.11(2)	5.86(9)	13.04(3)	9.66(5)	6.90(6)	6.86(7)	4.83(11)	3.11(14)	4.71(12)	6.24(8)
WBC	82.29(6)	89.87(3)	9.57(13)	92.27(1)	9.73(11)	73.56(8)	66.55(9)	54.00(10)	86.19(4)	6.38(12)	N/A(N/A)	78.67(7)	90.49(2)	
WDBC	75.46(5)	71.88(6)	14.93(13)	79.62(3)	50.52(9)	88.98(1)	43.72(10)	35.60(11)	84.78(2)	57.91(8)	6.57(14)	18.48(12)	66.11(7)	78.53(4)
Wilt	3.13(14)	3.62(12)	5.05(2)	3.64(1)	4.98(3)	3.84(9)	4.73(4)	5.53(1)	3.69(10)	4.14(7)	4.67(5)	4.00(8)	3.36(13)	4.23(6)
wine	30.87(4)	21.56(6)	7.77(13)	5.83(4)	8.45(10)	43.08(3)	8.43(11)	7.95(12)	45.71(2)	18.37(8)	21.14(7)	17.51(9)	48.82(1)	25.96(5)
WPBC	23.01(5)	22.93(6)	20.29(14)	21.32(12)	21.30(3)	23.04(4)	21.49(10)	25.37(3)	22.81(7)	21.38(11)	26.24(1)	22.49(8)	25.58(2)	22.42(9)
yeast	29.90(11)	29.84(12)	31.64(4)	30.93(7)	31.27(6)	32.75(3)	29.33(14)	29.96(3)	30.71(8)	31.36(5)	33.03(2)	29.92(10)	33.29(1)	29.80(13)
CIFAR10	10.59(6)	10.19(7)	13.02(1)	10.61(5)	11.61(2)	8.38(12)	11.13(3)	11.06(4)	8.77(11)	9.29(9)	8.05(13)	7.73(14)	9.72(8)	8.97(10)
FashionMNIST	31.42(6)	31.97(5)	16.85(13)	38.90(1)	20.73(11)	29.43(8)	33.87(2)	28.72(9)	30.32(7)	32.53(3)	17.43(2)	14.44(14)	27.32(10)	32.35(4)
MNIST-C	16.88(7)	17.72(6)	13.84(12)	27.62(1)	14.53(11)	15.46(10)	22.98(2)	15.68(9)	15.90(8)	18.24(4)	8.34(14)	11.37(13)	18.63(3)	17.99(5)
MVTec-AD	54.06(8)	51.44(10)	54.90(6)	58.52(1)	46.59(12)	55.22(5)	55.55(3)	51.48(9)	54.64(7)	55.44(4)	36.50(4)	45.66(13)	49.73(11)	56.04(2)
SVHN	8.66(5)	8.65(6)	9.24(2)	8.58(7)	8.97(3)	7.45(12)	9.46(1)	8.52(8)	7.61(11)	6.99(14)	7.29(13)	8.70(4)	8.10(9)	
Agnews	5.74(8)	5.69(9)	14.35(1)	7.02(5)	12.21(2)	5.58(10)	8.61(3)	8.40(4)	5.43(12)	4.45(14)	5.41(13)	5.93(7)	6.04(6)	
Amazon	5.85(9)	5.64(13)	5.72(11)	6.07(4)	5.74(10)	5.98(6)	6.23(2)	6.40(1)	6.08(3)	3.84(14)	5.65(12)	5.92(8)	5.95(7)	
Imdb	4.55(12)	4.44(14)	4.83(5)	4.75(6)	5.16(1)	4.74(8)	4.49(13)	4.70(9)	4.90(3)	5.06(2)	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)	9.20(2)	8.01(5)	7.98(6)	6.39(14)	6.72(13)	7.65(11)	7.88(7)
20news	7.97(5)	7.53(11)	9.13(1)	7.81(9)	9.02(2)	7.80(10)	8.54(3)	8.19(4)	7.95(6)	7.92(7)	7.82(8)	6.68(14)	7.37(12)	7.29(13)

Table D6: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 1\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRENet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTransformer	RF	LGB	XGB	CatB
ALOI	55.53(3)	59.13(2)	54.53(5)	47.03(15)	46.47(16)	55.06(4)	60.53(1)	49.31(13)	52.42(8)	48.50(14)	49.80(12)	51.74(9)	51.04(10)	50.60(11)	53.08(7)	53.22(6)
amnthyroid	75.67(9)	76.82(7)	72.20(11)	74.78(10)	75.95(8)	77.66(5)	92.89(2)	80.62(4)	59.25(15)	62.84(12)	52.04(16)	76.97(6)	69.93(14)	62.62(13)	89.91(3)	95.58(1)
backdoor	82.28(9)	91.98(3)	89.44(5)	92.91(2)	94.23(1)	82.60(7)	85.67(6)	63.10(13)	80.03(10)	89.69(4)	71.15(15)	62.89(14)	58.05(15)	37.26(16)	82.55(8)	76.47(11)
breastw	91.43(5)	88.19(3)	86.67(17)	74.21(10)	78.19(12)	80.76(6)	80.76(7)	82.37(14)	71.27(15)	52.52(15)	53.00(10)	95.67(2)	65.06(13)	96.04(4)	84.81(8)	97.40(1)
campaign	56.58(12)	64.32(7)	54.35(10)	61.49(11)	67.62(5)	58.45(11)	74.23(4)	47.30(14)	60.60(15)	61.68(16)	60.49(17)	61.68(8)	65.94(10)	58.45(14)	81.61(1)	81.61(1)
cardio	82.03(6)	69.21(12)	83.07(5)	89.74(2)	87.54(5)	79.57(8)	83.94(4)	58.89(5)	79.91(7)	61.55(14)	N/A(N/A)	75.38(9)	63.98(13)	73.56(11)	74.67(10)	94.01(1)
Cardiotocography	53.99(15)	64.94(12)	81.22(3)	81.91(2)	79.33(4)	70.67(10)	71.54(9)	76.00(7)	48.18(16)	65.16(11)	76.31(6)	60.95(13)	75.31(8)	78.81(5)	85.89(1)	
celeba	81.23(2)	54.65(14)	54.85(13)	73.84(5)	71.65(9)	71.59(7)	74.86(4)	50.71(15)	55.00(11)	62.62(8)	57.59(8)	55.63(10)	54.91(12)	29.53(16)	78.00(3)	81.50(1)
census	58.91(12)	60.87(11)	68.18(8)	73.80(8)	66.57(10)	67.16(9)	77.20(3)	50.66(16)	49.81(14)	76.76(4)	51.06(15)	70.22(7)	58.28(13)	71.26(6)	79.67(2)	84.53(1)
cover	42.98(16)	87.11(9)	98.60(4)	99.04(2)	98.63(3)	86.37(10)	92.92(8)	50.00(4)	80.88(11)	93.84(7)	62.00(13)	99.80(1)	62.41(12)	43.18(15)	96.09(6)	97.90(5)
donors	47.94(17)	97.54(4)	82.72(12)	99.71(2)	99.89(1)	96.82(6)	95.52(8)	95.34(9)	0.89(16)	70.69(12)	60.51(13)	97.89(3)	60.49(14)	77.58(11)	96.18(7)	96.98(5)
fault	63.85(2)	67.67(14)	63.79(8)	61.70(4)	63.19(6)	56.58(9)	60.30(10)	63.04(15)	47.09(14)	40.52(15)	53.25(9)	53.37(13)	54.21(12)	58.71(11)	60.27(14)	96.06(1)
fraud	90.52(5)	87.66(9)	92.32(2)	89.35(4)	91.88(1)	93.61(11)	93.98(6)	50.00(15)	92.44(1)	88.73(10)	75.13(13)	77.90(8)	77.90(8)	92.86(4)		
glass	67.58(15)	71.95(13)	85.97(6)	87.54(4)	90.77(2)	74.71(12)	83.81(11)	55.00(16)	83.87(9)	82.49(10)	86.30(5)	67.63(14)	84.94(7)	88.63(3)	91.09(1)	
Hepatitis	56.31(13)	62.67(11)	70.29(7)	68.49(9)	69.21(8)	50.98(15)	76.56(4)	53.62(14)	50.00(16)	74.39(6)	67.50(10)	64.53(12)	63.11(12)	82.29(2)	80.78(3)	83.14(1)
htpp	99.80(9)	99.88(8)	99.98(7)	100.00(1)	100.00(1)	99.99(6)	99.78(11)	81.67(14)	83.31(13)	0.10(16)	100.00(5)	83.33(12)	42.47(15)	100.00(1)	99.79(10)	
InternetAds	67.89(4)	71.41(2)	62.48(6)	51.66(3)	53.10(11)	60.17(7)	63.57(5)	50.28(4)	57.56(8)	42.13(15)	51.93(2)	N/A(N/A)	54.25(10)	55.15(9)	69.52(3)	77.18(1)
Ionosphere	91.98(1)	73.84(9)	77.56(3)	55.28(3)	54.48(14)	50.28(16)	75.87(5)	50.75(15)	74.09(8)	64.41(11)	65.63(10)	59.44(12)	74.91(6)	77.37(4)	87.27(2)	
landsat	45.19(16)	74.50(3)	57.12(14)	70.89(7)	73.67(4)	64.90(11)	75.46(2)	73.59(5)	56.48(15)	69.12(8)	59.53(13)	71.78(6)	60.09(12)	68.47(9)	67.12(10)	86.83(1)
letter	69.25(12)	64.47(11)	63.83(11)	50.13(4)	54.28(11)	65.36(14)	59.37(4)	48.23(14)	47.39(13)	60.77(14)	52.71(13)	53.54(13)	57.07(9)	59.71(7)	63.54(1)	86.78(1)
Lymphography	98.80(4)	80.40(9)	93.63(7)	80.65(9)	78.14(1)	82.39(8)	75.58(8)	56.77(14)	66.78(14)	50.00(16)	52.60(15)	47.79(12)	94.88(6)	72.04(13)	35.76(5)	99.65(1)
magic_gamma	52.62(14)	76.58(6)	73.47(7)	80.89(1)	80.46(2)	70.00(11)	76.70(5)	77.04(9)	50.00(16)	53.69(13)	51.56(15)	72.36(3)	59.53(12)	72.14(9)	70.33(10)	72.51(8)
mammography	77.28(8)	84.40(6)	88.92(3)	85.87(4)	84.47(5)	89.16(2)	78.37(7)	73.34(10)	60.78(13)	19.48(6)	61.43(12)	72.23(9)	58.85(14)	39.41(5)	66.95(1)	89.83(1)
mnist	69.68(11)	75.32(8)	79.96(5)	76.71(7)	71.04(10)	88.35(2)	64.07(14)	69.29(12)	85.93(7)	63.76(13)	74.27(9)	64.05(15)	61.14(16)	80.02(4)	91.19(1)	
mask	99.12(3)	86.41(11)	84.71(12)	88.19(10)	88.26(9)	89.38(8)	50.00(15)	43.08(16)	95.16(6)	92.35(3)	97.39(4)	63.78(13)	55.51(4)	89.42(7)	99.44(2)	
optdigits	45.68(15)	84.32(10)	99.23(4)	94.15(8)	94.15(8)	50.74(14)	87.83(4)	82.84(11)	67.12(3)	98.55(5)	75.39(12)	43.88(6)	94.92(7)	96.46(6)		
PageBlocks	47.79(6)	88.76(4)	88.50(2)	70.47(8)	73.14(7)	54.12(13)	82.95(5)	64.43(11)	35.02(16)	36.07(16)	70.52(20)	55.59(13)	65.55(10)	86.54(3)	93.80(1)	
pedigree	82.62(12)	82.47(11)	82.63(11)	80.67(12)	82.38(9)	82.63(11)	82.38(9)	83.72(13)	84.14(14)	79.95(12)	82.75(15)	86.63(15)	81.64(11)	87.20(14)	88.57(13)	94.21(1)
Pima	59.71(8)	58.30(10)	67.54(2)	59.98(7)	58.30(11)	58.78(8)	64.07(3)	26.31(12)	48.99(15)	23.38(16)	54.04(8)	63.06(4)	54.29(13)	62.21(5)	71.70(6)	71.28(1)
satellite	68.69(13)	78.70(6)	75.63(2)	82.60(1)	81.44(2)	77.99(7)	80.29(4)	78.99(5)	55.48(16)	75.44(9)	60.82(12)	91.64(8)	59.84(12)	42.36(6)	94.80(5)	98.57(1)
satimage-2	96.55(7)	95.66(9)	96.40(8)	98.37(3)	97.99(4)	98.44(2)	97.88(5)	50.00(14)	47.99(15)	79.10(12)	94.53(11)	96.78(6)	69.75(13)	49.08(16)	94.64(10)	98.77(1)
shuttle	73.14(14)	88.92(1)	88.92(1)	76.28(1)	87.80(6)	89.13(5)	88.74(2)	72.06(12)	70.28(13)	89.43(1)	80.61(16)	87.75(14)	81.73(12)	84.72(11)	87.43(10)	95.43(1)
skin	52.14(16)	92.43(2)	92.43(2)	92.43(1)	92.43(2)	92.43(1)	92.43(1)	53.39(1)	53.30(7)	37.51(7)	53.30(7)	53.30(7)	53.30(7)	37.51(7)	53.30(7)	92.43(1)
SpamBase	65.24(14)	64.94(13)	64.68(14)	99.98(2)	99.98(2)	60.11(14)	99.23(8)	60.00(15)	52.98(13)	52.98(13)	52.98(13)	52.98(13)	52.98(13)	52.98(13)	52.98(13)	99.98(2)
speech	41.22(7)	57.34(9)	94.19(5)	94.19(5)	94.19(5)	94.19(5)	94.19(5)	53.02(11)	53.02(11)	53.02(11)	53.02(11)	53.02(11)	53.02(11)	53.02(11)	53.02(11)	94.19(5)
stamps	73.48(11)	74.48(10)	93.99(7)	97.78(6)	80.39(8)	53.60(13)	61.84(14)	0.37(16)	100.00(1)	99.70(5)	47.10(16)	87.58(6)	63.60(12)	58.72(14)	83.45(5)	87.70(1)
InternetAds	46.05(3)	42.68(4)	41.92(5)	33.91(10)	34.32(9)	39.10(6)	36.29(7)	38.18(2)	34.55(8)	17.30(10)	36.29(7)	34.55(8)	34.55(8)	34.55(8)	34.55(8)	35.84(8)
Ionosphere	89.55(1)	69.44(3)	64.46(8)	47.03(11)	53.39(7)	49.57(14)	53.38(6)	46.42(16)	63.89(9)	68.45(4)	57.52(11)	57.75(5)	48.79(15)	65.26(7)	52.77(11)	80.00(1)
landsat	11.72(16)	3.79(13)	2.94(15)	2.94(15)	14.42(7)	14.42(7)	11.64(7)	2.64(16)	7.71(10)	7.81(10)	8.67(11)	5.42(16)	5.42(16)	3.10(14)	14.31(4)	12.62(5)
letter	73.58(11)	6.78(14)	9.57(11)	9.57(11)	1.26(14)	27.98(10)	34.32(7)	61.61(14)	1.04(16)	9.94(11)	22.70(14)	12.70(14)	21.29(12)	20.41(14)	30.37(12)	70.68(1)
Lymphography	56.50(6)	39.25(9)	31.33(9)	14.25(7)	14.25(7)	14.87(5)	14.36(6)	18.55(2)	6.25(16)	34.15(7)	17.08(14)	18.88(14)	12.99(16)	13.46(14)	11.44(12)	77.76(1)
magic_gamma	39.96(16)	64.64(5)	66.45(4)	71.97(2)	72.46(1)	62.48(20)	78.02(3)	57.51(4)	53.52(14)	48.51(12)	45.12(15)	45.12(15)	45.12(15)	45.12(15)	45.12(15)	67.37(10)
mammography	15.75(13)	32.24(7)	40.20(4)	40.49(3)	36.69(6)	53.56(1)	23.52(10)	29.90(8)	18.01(12)	1.37(16)	22.86(13)	42.42(4)	23.30(14)	23.44(15)	37.72(5)	47.30(6)
mnist	17.88(16)	32.84(10)	48.06(7)	56.31(2)	56.31(5)	50.16(6)	51.69(5)	22.77(15)	35.17(16)	45.26(15)	35.17(16)	35.17(16)	35.17(16)	35.17(16)	35.17(16)	35.84(5)
optdigits	2.83(16)	27.71(13)	76.11(13)	76.11(13)	76.11(13)	45.27(17)	77.76(13)	84.02(16)	34.22(16)	57.20(13)	77.20(13)	82.30(14)	56.91(15)	56.55(14)	49.15(14)	50.55(1)
PageBlocks	41.22(7)	57.16(3)	57.34(2)	49.56(18)	49.56(18)	33.40(22)	47.41(14)	34.78(6)	1.13(16)	11.66(17)	47.41(14)	33.80(16)	33.80(16)	23.02(15)	33.80(11)	57.14(4)
Pima	43.77(11)	44.26(10)	50.96(2)	45.76(7)	44.44(2)	42.87(12)	50.02(3)	44.93(18)	33.26(13)	23.16(11)	41.13(13)	47.15(6)	38.80(14)	47.18(5)	47.40(4)	53.50(1)
satellite	65.05(10)	67.43(9)	71.23(5)	74.41(1)	74.41(2)	75.32(5)	70.55(6)	73.20(4)	44.82(16)	60.43(12)	50.63(14)	70.54(7)	48.47(15)	60.63(11)	70.24(8)	
satimage-2	43.15(10)	40.06(11)	85.70(7)	89.93(9)	92.20(2)	50.83(10)	85.52(7)	90.40(10)	88.60(12)	26.21(15)	25.15(13)	88.44(11)	20.75(16)	54.32(13)	91.09(8)	94.95(8)
shuttle	46.19(14)	96.40(4)	97.21(1)	96.52(2)	96.50(3)	93.39(5)	50.16(6)	51.69(5)	34.05(15)	31.35(16)	31.35(16)	31.35(16)	31.35(16)	31.35(16)	31.35(16)	92.77(7)
skin	21.90(16)	83.45(14)	61.61(12)	60.85(10)	60.85(10)	81.73(11)	84.85(11)	84.30(15)	40.98(16)	31.03(15)	31.					

Table D8: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 5\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in bold.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRENet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTransformer	RF	LGB	XGB	CatB
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)	49.80(16)	71.31(15)	77.35(3)	95.13(4)	82.25(8)	93.14(5)	95.24(3)	98.51(1)
backdoor	81.81(12)	95.13(3)	89.33(10)	95.58(2)	94.37(4)	96.80(1)	91.07(8)	79.06(13)	86.80(11)	93.12(7)	66.31(16)	89.84(9)	76.68(14)	76.63(15)	93.34(6)	94.34(5)
breastwt	92.28(11)	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(3)	85.59(4)	97.53(5)	95.12(8)	99.10(3)
campagne	64.43(14)	69.71(10)	57.88(6)	81.03(5)	80.65(6)	71.01(8)	88.41(2)	69.02(11)	68.58(12)	70.18(9)	59.17(15)	67.08(13)	72.78(7)	84.34(4)	86.48(3)	88.89(1)
cardio	83.20(10)	87.82(4)	80.13(7)	95.90(1)	90.98(2)	88.79(5)	92.72(3)	71.04(12)	70.84(12)	70.71(13)	88.75(14)	88.03(5)	84.98(4)	89.97(3)	96.12(2)	96.83(2)
Cardiotocography	54.77(11)	81.24(12)	91.89(3)	91.17(1)	80.58(7)	91.49(10)	83.73(8)	84.05(12)	84.67(10)	60.44(15)	81.73(11)	79.49(14)	83.03(9)	84.55(7)	89.12(4)	89.12(4)
celeba	68.79(12)	75.40(10)	56.47(11)	91.01(1)	90.64(3)	68.63(6)	87.45(5)	74.33(11)	66.62(5)	90.93(2)	58.14(14)	69.28(5)	63.12(13)	83.72(9)	83.90(8)	84.69(7)
census	59.66(14)	69.05(13)	69.20(12)	81.03(5)	75.34(10)	76.32(8)	87.48(2)	57.28(15)	75.94(6)	75.94(9)	56.08(6)	78.81(7)	70.47(11)	83.37(4)	88.04(1)	86.81(3)
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)	89.40(12)	96.38(8)	68.67(7)	99.92(1)	84.32(13)	57.48(15)	94.08(10)	99.24(5)
donors	68.62(15)	99.92(4)	82.81(3)	99.99(1)	99.95(3)	99.99(2)	99.59(7)	99.42(8)	99.62(6)	95.70(11)	99.90(5)	77.22(14)	95.13(12)	97.90(10)	99.29(9)	99.29(9)
fault	64.26(12)	88.76(4)	67.54(2)	67.54(2)	67.54(2)	67.54(2)	67.54(2)	63.65(5)	62.15(7)	64.00(10)	50.20(16)	65.99(7)	65.07(8)	66.15(6)	70.15(7)	72.03(1)
frad	90.53(14)	87.89(9)	93.18(2)	88.77(7)	87.60(10)	87.53(11)	89.84(6)	50.94(15)	68.34(11)	49.57(10)	59.39(14)	88.38(3)	84.92(11)	44.99(16)	87.31(9)	90.19(4)
glass	67.84(16)	84.34(11)	82.27(11)	91.37(9)	94.46(4)	93.72(5)	95.73(3)	92.73(6)	80.63(14)	82.59(12)	89.96(10)	71.48(8)	78.85(5)	92.17(7)	95.97(2)	98.37(1)
Hepatitis	56.94(15)	74.83(12)	88.10(4)	86.62(5)	83.71(8)	72.08(13)	85.44(6)	77.79(11)	50.00(16)	85.11(7)	80.48(9)	82.99(9)	69.71(14)	89.29(3)	89.66(2)	91.72(1)
http	99.80(11)	99.90(10)	99.98(8)	100.00(1)	100.00(1)	99.99(7)	98.33(12)	98.33(12)	83.31(14)	0.14(16)	100.00(1)	100.00(1)	100.00(6)	39.87(15)	100.00(1)	99.97(9)
InternetAds	67.94(8)	72.98(4)	79.79(2)	61.68(11)	57.82(2)	63.57(10)	75.77(3)	51.14(9)	62.65(5)	74.33(11)	56.12(5)	74.13(14)	69.28(5)	63.12(13)	68.68(6)	67.87(9)
ionosphere	81.80(9)	94.65(4)	68.63(7)	65.63(11)	63.46(15)	67.53(14)	62.45(2)	70.27(13)	80.25(9)	76.28(8)	66.98(7)	84.70(5)	82.67(6)	88.44(1)	88.44(1)	95.12(1)
landesk	46.33(16)	55.54(9)	59.94(15)	70.39(12)	52.21(12)	87.34(2)	87.34(2)	71.02(13)	70.27(10)	50.00(15)	74.07(14)	74.33(11)	76.94(10)	76.59(11)	84.53(7)	84.53(7)
letter	69.64(4)	73.09(3)	66.43(5)	65.01(7)	63.18(1)	64.18(8)	83.13(1)	48.83(14)	40.02(16)	42.04(15)	76.96(2)	59.82(3)	62.33(12)	65.92(6)	63.80(9)	63.50(10)
Lymphography	96.84(2)	86.90(8)	92.94(4)	80.25(9)	79.68(10)	93.47(3)	77.40(13)	73.19(14)	50.00(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(15)	82.08(6)	78.24(10)	82.76(6)	83.38(4)	84.04(3)	84.58(2)	76.10(13)	49.30(16)	85.08(1)	70.15(8)	81.91(7)	78.06(12)	79.11(9)	81.85(8)	81.85(8)
mammography	75.46(13)	90.90(6)	91.90(3)	92.92(0)	93.59(3)	94.77(11)	91.28(4)	72.22(10)	76.27(9)	50.20(15)	72.32(11)	67.59(15)	82.84(6)	86.55(11)	92.93(2)	97.02(1)
mnist	97.03(16)	98.91(12)	98.91(12)	98.91(12)	98.91(12)	98.91(12)	98.91(12)	99.00(12)	98.92(12)	98.92(12)	98.92(12)	98.92(12)	98.92(12)	98.92(12)	98.92(12)	98.92(12)
optdigits	46.44(16)	94.98(10)	99.67(4)	99.98(1)	99.94(3)	99.94(3)	98.23(8)	80.46(13)	90.94(12)	99.48(5)	75.71(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(7)	90.38(7)	89.67(9)	94.73(2)	88.46(10)	58.16(16)	90.15(8)	82.94(3)	88.03(1)	76.84(14)	92.74(5)	94.28(3)	96.80(1)
pima	56.82(12)	97.41(11)	99.75(2)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)	99.76(1)
pendigits	49.53(16)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)	93.30(7)
satellite	60.11(13)	63.23(12)	73.18(4)	78.02(1)	76.63(4)	63.55(11)	69.92(8)	71.13(6)	47.99(5)	32.37(16)	51.59(4)	68.69(9)	67.36(10)	71.80(5)	73.80(3)	73.80(3)
satellite2	69.49(15)	88.05(4)	81.31(1)	84.34(10)	84.41(9)	85.83(7)	85.83(7)	79.37(13)	61.94(16)	87.79(5)	87.82(4)	82.71(11)	86.74(12)	85.20(8)	90.22(1)	89.23(1)
satellite-g	74.64(15)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)	88.41(1)
shuttle	77.41(16)	98.71(4)	98.83(3)	97.57(0)	97.68(7)	99.73(1)	94.63(14)	97.04(13)	97.49(9)	97.44(10)	78.01(15)	97.16(12)	98.00(5)	98.85(2)	98.85(2)	98.85(2)
skin	52.78(16)	99.44(14)	99.56(15)	95.72(16)	95.24(12)	98.56(8)	99.22(5)	94.89(13)	99.60(3)	97.23(10)	99.83(1)	99.06(6)	78.23(15)	99.00(7)	98.53(9)	99.72(2)
smt	51.69(14)	82.74(3)	76.47(8)	75.99(9)	75.24(10)	74.57(11)	70.87(5)	62.98(12)	86.40(1)	82.83(13)	79.07(6)	83.32(2)	81.43(16)	78.85(7)	78.85(7)	78.85(7)
SpamBase	53.66(16)	70.72(13)	83.55(8)	72.25(3)	80.28(11)	82.57(12)	83.18(1)	78.45(11)	73.90(4)	51.11(14)	65.29(6)	45.15(5)	63.29(9)	61.57(11)	64.77(7)	70.42(5)
speech	47.49(15)	51.11(11)	51.13(10)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)	51.11(11)
spambase	52.57(14)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)	50.98(11)
thyroid	92.65(11)	95.21(10)	99.55(1)	99.53(4)	99.52(5)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.53(1)	99.61(1)
vertebral	56.18(16)	68.41(4)	72.26(6)	71.58(7)	71.29(8)	62.91(11)	76.56(4)	51.90(9)	61.29(12)	55.90(15)	14.35(16)	60.28(8)	73.35(11)	77.35(1)	63.15(13)	77.34(1)
vowels	78.54(16)	88.47(6)	89.12(6)	92.49(1)	86.69(5)	87.82(3)	73.29(12)	55.90(15)	46.15(16)	76.82(10)	78.62(13)	77.98(4)	81.68(8)	81.23(9)	76.14(11)	86.61(1)
Waveform	52.99(16)	68.05(10)	80.58(5)	84.84(1)	82.30(3)	64.09(13)	74.68(7)	77.42(6)	60.62(5)	61.03(14)	72.60(8)	80.85(4)	64.74(12)	61.81(11)	71.76(9)	84.82(2)
WBC	92.78(16)	94.96(4)	97.88(5)	98.38(2)	98.28(3)	98.69(4)	98.44(3)	98.72(1)	98.22(1)	98.83(1)	98.63(1)	98.79(2)	98.71(1)	98.80(1)	98.81(1)	98.81(1)
WDBC	77.57(14)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)	95.96(1)
WIS	74.31(16)	78.34(10)	80.89(3)	99.00(1)	99.12(3)	74.22(8)	45.46(13)	45.98(12)	94.68(5)	74.33(12)	27.45(13)	58.26(6)	45.13(14)	53.88(15)	53.88(16)	54.10(1)
wine	14.98(14)	37.92(13)	99.15(3)	100.00(1)	98.44(4)	78.20(11)	88.68(7)	81.11(5)	4.76(16)	80.39(9)	98.23(5)	51.74(12)	78.46(10)	82.05(8)	93.48(6)	93.48(6)
Wdbc	22.58(15)	29.86(13)	30.72(12)	40.76(8)	45.11(4)	37.69(10)	45.79(3)	30.99(11)	22.58(16)	25.35(14)	51.48(1)	42.56(7)	39.91(9)	45.02(5)	47.26(2)	44.02(6)
yeast	33.03(16)	36.31(13)	34.39(14)	46.72(1)	46.60(3)	45.30(5)	40.33(10)	46.69(2)	45.58(4)	34.25(15)	32.77(7)	40.52(9)	41.30(8)	39.80(11)	45.11(6)	45.11(6)
CIFFAR10	92.00(14)	12.26(12)	19.48(6)	21.16(1)	20.79(3)	20.18(6)	20.99(2)	5.20(16)	13.12(10)	15.32(9)	12.47(11)	9.65(13)	7.28(15)	15.73(8)	20.75(4)	18.09(7)
FashionMNIST	26.10(10)	43.27(4)	52.34(5)	60.74(1)	62.12(2)	53.73(6)	57.74(3)	7.35(16)	26.13(13)	34.78(11)	29.41(12)	32.22(11)	7.74(13)	41.94(11)	55.06(4)	53.05(5)
MNISTC	23.82(11)	50.04(10)	1.89(1)	9.16(1)	32.34(1)	62.12(3)	64.81(10)	52.99(7)	48.40(9)	51.65(9)	49.38(11)	49.38(11)	49.38(11)	49.38(11)	49.38(11)	49.38(11)
MVTC-AD	57.54(5)	52.50(11)	54.55(8)	56.29(6)	56.07(5)	52.41(12)	61.21(3)	28.19(16)	36.03(5)	50.02(

Table D10: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 10\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRENet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTransformer	RF	LGB	XGB	CatB	
ALOI	56.59(7)	56.63(6)	54.54(8)	49.55(13)	49.50(14)	53.11(9)	72.10(1)	47.78(15)	47.05(16)	51.34(10)	49.79(12)	50.39(11)	60.67(2)	57.98(4)	58.95(3)	56.85(5)	
amnthyroid	82.83(1)	87.41(9)	82.05(4)	83.31(10)	82.78(2)	93.07(7)	98.21(3)	82.54(13)	40.43(16)	80.53(15)	98.88(8)	98.67(2)	93.72(6)	97.39(5)	97.56(4)	98.85(1)	
backdoor	83.69(14)	96.47(5)	89.40(12)	97.84(1)	96.39(6)	96.73(3)	96.17(7)	79.46(15)	94.03(9)	95.82(8)	58.20(16)	93.93(10)	87.71(13)	93.93(11)	97.22(2)	96.69(4)	
breastwt	92.59(13)	94.61(11)	97.92(8)	99.49(3)	98.79(5)	88.12(4)	97.88(9)	99.75(1)	79.54(15)	99.54(2)	60.44(16)	98.18(7)	93.69(12)	99.01(4)	97.43(10)	98.78(6)	
campagne	73.02(10)	71.02(13)	57.92(16)	83.85(5)	83.84(6)	72.49(12)	90.77(1)	75.62(20)	68.88(14)	72.51(11)	62.26(15)	79.61(8)	80.15(7)	88.25(6)	88.72(3)	90.08(2)	
carbs	82.86(1)	86.45(5)	90.37(6)	98.86(5)	98.90(7)	87.37(7)	97.96(1)	90.92(15)	86.97(16)	87.99(14)	80.20(16)	86.96(10)	91.29(14)	92.10(19)	94.21(1)	98.84(1)	
Cardiotocography	55.15(2)	86.65(10)	18.80(4)	33.15(2)	85.54(1)	85.53(3)	89.13(7)	89.77(9)	88.85(6)	83.19(14)	91.02(5)	86.34(5)	85.64(12)	87.23(8)	87.03(9)	91.85(3)	
celeba	72.16(13)	80.00(11)	57.53(11)	93.69(2)	93.63(5)	88.66(9)	92.90(5)	88.89(8)	79.76(1)	94.03(1)	62.02(15)	93.63(4)	71.76(3)	91.10(14)	90.03(6)	89.17(7)	87.29(10)
census	60.45(16)	71.90(12)	69.07(13)	84.88(5)	79.66(6)	79.20(7)	88.53(1)	62.38(14)	79.02(10)	76.73(11)	61.42(16)	77.63(11)	79.14(8)	85.87(4)	86.07(3)	87.91(2)	
cover	45.42(16)	96.71(10)	99.92(2)	99.88(3)	99.87(4)	95.43(1)	97.76(9)	98.72(6)	92.73(12)	97.86(8)	73.39(14)	99.92(1)	92.37(13)	61.22(15)	98.25(7)	99.50(5)	
donors	69.68(16)	99.93(2)	82.79(5)	99.99(1)	99.95(5)	99.98(4)	99.84(6)	99.42(11)	99.72(8)	99.84(6)	99.43(10)	99.99(10)	90.19(14)	99.01(13)	99.28(12)	99.72(9)	
fault	64.45(14)	72.45(7)	75.19(4)	82.16(1)	77.65(2)	72.05(6)	73.81(2)	67.35(16)	63.85(15)	73.92(2)	55.92(6)	71.13(9)	73.51(4)	84.01(10)	84.23(16)	84.41(12)	
gas	90.53(3)	87.00(7)	92.22(1)	87.01(0)	85.45(5)	80.95(5)	81.86(1)	90.25(5)	89.98(6)	89.29(6)	90.19(6)	89.20(6)	89.20(10)	89.20(8)	89.20(2)	97.79(1)	
Hepatitis	65.34(15)	85.24(13)	93.02(8)	95.68(3)	93.86(6)	88.61(12)	94.25(5)	91.89(10)	50.06(16)	92.91(9)	91.54(11)	93.12(7)	85.17(14)	96.29(2)	95.32(4)	96.35(1)	
http	99.80(11)	99.93(9)	99.98(7)	100.00(1)	100.00(1)	99.84(10)	98.33(12)	98.33(12)	83.31(14)	0.17(16)	100.00(1)	100.00(6)	44.39(15)	100.00(1)	99.97(8)	100.00(1)	
InternetAds	68.04(11)	75.03(7)	82.43(2)	70.71(9)	66.90(2)	70.04(10)	76.69(3)	53.78(5)	75.03(8)	66.15(13)	57.26(6)	70.46(16)	75.62(6)	76.23(5)	76.41(4)	84.26(1)	
ionosphere	89.20(5)	85.79(5)	94.55(1)	78.35(1)	78.35(1)	78.35(1)	95.73(1)	95.73(1)	77.73(1)	82.70(15)	82.70(15)	87.08(15)	83.13(16)	93.40(15)	92.34(16)	94.42(2)	
lansup	47.49(16)	51.43(1)	57.89(7)	68.71(2)	70.28(1)	59.16(5)	57.27(9)	58.63(6)	57.79(6)	52.67(12)	52.58(13)	68.85(5)	51.69(5)	52.34(4)	56.02(10)	63.20(4)	
pendigits	57.10(16)	99.38(9)	99.78(4)	99.86(3)	99.97(1)	99.97(1)	99.97(1)	99.54(7)	89.42(14)	99.30(10)	97.21(3)	99.91(1)	90.75(12)	85.90(15)	99.47(8)	99.73(5)	
Pima	60.29(13)	66.07(12)	69.69(9)	76.24(2)	77.35(1)	71.49(8)	73.83(6)	74.76(5)	48.36(16)	58.11(14)	49.26(15)	69.65(10)	75.29(4)	73.70(7)	75.53(3)	80.00(1)	
satimage-2	71.51(36)	91.07(4)	80.74(1)	84.88(10)	83.94(7)	86.78(9)	91.02(4)	79.02(14)	76.90(15)	81.87(12)	92.42(1)	89.86(16)	89.62(26)	88.53(16)	92.46(2)	92.61(2)	
shuttle	96.34(8)	97.99(4)	99.02(5)	99.56(1)	99.55(1)	99.42(1)	99.70(1)	94.00(14)	93.67(15)	96.01(14)	96.90(16)	97.13(5)	85.98(14)	98.69(3)	93.40(8)	99.72(1)	
skin	63.73(16)	99.01(3)	98.72(5)	97.55(10)	97.69(8)	98.23(7)	99.88(1)	93.34(1)	87.58(12)	97.47(11)	97.64(9)	87.27(12)	90.47(14)	89.67(6)	98.94(4)	99.51(2)	
smt	52.26(16)	99.64(3)	90.29(1)	95.70(11)	95.30(2)	92.22(10)	99.58(4)	94.81(13)	99.18(13)	99.53(14)	92.46(6)	93.15(1)	95.35(4)	96.66(6)	98.09(2)	99.84(1)	
SpamBase	54.01(16)	80.28(13)	87.14(9)	91.35(6)	92.96(5)	89.84(3)	88.33(9)	84.94(1)	85.20(15)	62.98(12)	86.40(1)	52.83(12)	79.07(6)	83.32(2)	43.01(16)	81.12(4)	
speech	47.49(16)	54.23(11)	57.89(7)	68.71(2)	70.28(1)	59.16(5)	57.27(9)	58.63(6)	57.79(6)	52.67(12)	52.58(13)	68.85(5)	51.69(5)	52.34(4)	56.02(10)	63.20(4)	
spells	71.38(15)	85.80(11)	94.88(1)	94.84(1)	94.84(1)	94.84(1)	94.84(1)	87.72(15)	89.84(15)	89.84(15)	87.72(15)	89.84(15)	87.72(15)	89.84(15)	89.84(15)	95.60(1)	
thyroid	92.77(16)	97.10(10)	99.70(2)	99.73(1)	99.51(2)	99.49(3)	99.49(4)	99.49(4)	99.41(10)	99.41(14)	99.41(14)	99.41(14)	99.41(14)	99.41(14)	99.41(14)	99.50(1)	
vertebral	32.77(16)	68.06(11)	56.21(15)	76.36(9)	78.26(7)	72.74(8)	85.83(2)	83.34(4)	78.55(10)	67.61(12)	62.55(16)	80.83(26)	65.78(13)	82.36(5)	87.31(1)	84.73(3)	
vowels	78.52(12)	82.64(12)	92.68(5)	95.69(3)	96.99(1)	93.00(4)	95.74(2)	88.52(9)	79.63(14)	15.44(16)	81.05(12)	84.07(11)	88.89(8)	91.02(7)	92.56(6)	99.84(1)	
Waveform	52.95(10)	74.37(13)	83.17(5)	88.17(3)	88.84(2)	77.71(12)	83.77(4)	81.92(7)	80.17(10)	62.68(15)	78.54(11)	80.45(9)	70.38(14)	81.63(8)	82.52(6)	89.82(1)	
WBC	93.57(9)	98.18(13)	98.82(4)	98.82(1)	98.82(1)	98.82(1)	98.82(1)	97.73(12)	97.73(12)	98.82(1)	98.82(1)	98.82(1)	98.82(1)	98.82(1)	98.82(1)	99.37(1)	
WDBC	47.57(16)	95.95(12)	95.10(3)	100.00(1)	99.95(1)	99.96(1)	99.97(1)	99.97(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	100.00(1)	
Wii	47.30(16)	50.92(5)	53.39(5)	59.08(2)	60.14(1)	59.99(2)	60.20(2)	59.81(12)	58.63(11)	58.63(14)	59.98(14)	58.63(14)	62.05(19)	89.13(6)	89.59(6)	95.66(1)	
wine	66.14(15)	89.42(12)	99.87(5)	100.00(1)	100.00(1)	86.04(8)	84.08(1)	89.19(4)	87.23(12)	80.05(16)	86.16(16)	87.57(7)	90.36(10)	77.25(12)	92.98(10)	91.78(11)	
WPBC	47.62(15)	64.81(12)	63.93(1)	77.76(1)	77.43(3)	73.90(8)	75.00(6)	62.06(14)	75.15(4)	72.46(9)	66.25(15)	72.15(1)	69.37(10)	74.17(7)	77.52(2)	77.39(4)	
yeast	48.58(16)	59.53(16)	47.25(1)	64.84(9)	65.55(3)	63.30(6)	59.36(3)	65.71(2)	59.31(13)	53.31(14)	53.81(13)	46.61(5)	57.64(12)	60.25(7)	59.54(10)	66.23(1)	
CIFAR10	78.76(9)	94.49(3)	94.46(3)	94.54(3)	94.54(3)	94.66(16)	94.66(16)	86.05(6)	87.45(12)	87.45(12)	86.05(6)	87.45(12)	87.45(12)	87.45(12)	87.45(12)	92.77(7)	
FashionMNIST	60.39(13)	65.50(11)	73.80(6)	74.61(4)	76.98(7)	72.96(6)	76.66(1)	53.15(16)	53.77(11)	53.77(11)	53.77(11)	53.77(11)	53.77(11)	53.77(11)	53.77(11)	60.26(11)	
MNISTC	75.79(15)	88.62(10)	91.59(2)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	92.23(1)	94.95(1)	
MVTe-AD	74.28(8)	73.68(9)	73.33(1)	80.21(5)	76.70(1)	76.70(2)	82.61(2)	56.36(16)	68.41(14)	71.60(11)	61.40(15)	70.14(13)	74.67(7)	80.55(4)	82.80(3)	83.92(1)	
SVHN	57.61(14)	63.11(16)	70.74(2)	74.78(1)	74.67(4)	70.35(4)	70.22(3)	50.09(16)	63.73(10)	62.61(12)	61.79(14)	63.78(9)	54.21(15)	68.17(6)	70.65(3)	76.61(3)	
Agnews	58.77(14)	74.98(12)	88.09(3)	88.18(2)	88.32(7)	84.99(5)	82.23(15)	88.23(18)	84.17(18)	86.98(13)	92.04(12)	87.98(12)	87.98(12)	84.23(16)	84.34(16)	81.34(9)	
Amazon	57.80(14)	59.22(13)	77.98(5)	81.80(9)	81.80(9)	81.62(1)	81.79(7)	72.01(13)	57.46(15)	67.46(15)	72.01(13)	72.01(13)	81.46(16)	70.63(16)	80.03(10)	88.60(9)	
Imdb	49.39(16)	63.87(12)	74.92(7)	80.80(2)	80.80(3)	76.39(5)	74.83(8)	55.83(4)	74.81(14)	61.92(14)	88.77(3)	60.42(15)	78.73(7)	79.60(11)	86.29(6)	90.76(3)	
Yelp	66.76(13)	68.39(12)	87.34(4)	89.48(2)	89.03(5)	85.52(5)	82.23(7)	53.76(13)	81.09(8)	90.32(1)	66.67(14)	74.45(1)	55.79(15)	77.55(10)	82.63(6)	80.41(9)	
20news	57.18(14)	64.91(9)	57.15(15)	72.45(1)	69.66(6)	70.61(10)	65.95(8)	65.83(8)	58.24(13)	65.84(16)	65.84(16)	65.84(16)	65.84(16)	65.84(16)	65.84(16)	65.84(16)	

Table D11: AUCPR of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 10\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRENet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTransformer	RF	LGB	XGB	CatB
ALOI	3.92(12)	6.36(5)	5.12(9)	4.65(10)	5.30(7)	3.70(14)	8.15(1)									

Table D12: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 25\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRENet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FTTransformer	RF	LGB	XGB	CatB		
ALOI	55.25(8)	59.20(6)	53.09(11)	51.65(3)	51.44(4)	57.54(7)	75.96(1)	49.27(15)	53.86(9)	52.16(12)	47.74(16)	53.45(10)	71.83(2)	65.23(3)	64.77(4)	62.78(5)		
amnthyroid	80.96(15)	93.06(9)	82.78(2)	82.39(3)	82.11(4)	96.95(7)	98.78(4)	82.95(11)	50.73(16)	91.29(10)	96.98(8)	99.05(1)	98.58(6)	98.70(5)	98.83(3)	98.97(2)		
backdoor	82.96(15)	97.07(8)	89.49(13)	97.49(6)	96.45(9)	98.00(4)	97.20(7)	76.77(11)	94.51(12)	95.98(10)	86.10(14)	94.98(11)	97.92(5)	98.40(2)	98.39(3)	98.92(1)		
breastwt	93.13(15)	97.44(12)	98.72(8)	99.67(7)	99.27(5)	93.83(14)	99.14(6)	99.75(1)	94.57(13)	99.42(3)	80.40(16)	98.45(10)	98.23(11)	99.30(4)	98.72(9)	99.11(7)		
campaign	66.57(15)	73.12(14)	57.92(16)	85.83(7)	85.95(6)	76.42(11)	92.20(1)	79.15(19)	73.59(13)	77.29(10)	75.45(12)	84.62(8)	87.85(5)	90.41(6)	89.16(4)	91.23(2)		
car	3.83(15)	95.25(14)	96.44(10)	98.00(9)	96.29(11)	97.48(6)	96.23(5)	93.15(13)	86.79(15)	89.96(14)	94.26(10)	96.30(9)	95.92(1)	95.96(2)	96.33(3)	96.92(2)		
Cardiotocography	56.46(11)	92.52(12)	84.43(6)	94.03(5)	94.20(4)	90.40(3)	93.81(16)	93.20(10)	87.28(12)	89.08(4)	95.25(9)	94.05(8)	93.34(9)	93.50(8)	92.89(11)	95.13(1)		
celleba	76.18(15)	89.65(11)	57.39(9)	94.91(1)	95.17(3)	90.65(9)	89.68(10)	85.43(12)	94.44(6)	77.31(4)	95.20(2)	81.65(13)	94.51(5)	91.86(8)	92.40(7)			
census	46.38(16)	78.77(12)	69.41(9)	90.33(3)	88.07(6)	82.45(10)	91.07(1)	70.09(13)	83.98(9)	86.41(8)	98.04(11)	82.44(7)	89.95(4)	89.72(5)	90.85(2)			
cover	42.98(16)	99.86(5)	99.94(1)	99.92(2)	99.94(2)	99.83(7)	99.68(12)	99.75(11)	98.97(13)	99.77(9)	96.56(14)	99.81(8)	99.85(6)	90.19(15)	99.76(10)	99.91(4)		
donors	51.09(16)	100.00(1)	82.84(1)	99.99(5)	99.94(4)	99.98(3)	99.93(6)	99.42(14)	99.57(12)	99.89(8)	99.43(3)	99.93(3)	99.92(7)	99.58(1)	99.88(9)	99.82(10)	99.93(5)	
fault	65.21(16)	75.68(13)	70.79(10)	74.19(7)	72.33(6)	70.46(4)	71.76(1)	77.71(13)	71.15(13)	74.87(4)	77.91(2)	66.78(15)	75.57(7)	76.47(5)	76.18(6)	80.10(1)		
fraud	91.06(3)	90.21(6)	92.38(1)	92.64(7)	90.47(6)	79.07(10)	90.46(13)	76.83(10)	90.68(12)	90.48(12)	83.32(11)	88.20(8)	92.71(16)	85.83(1)	90.20(2)			
glass	68.06(16)	91.67(12)	86.42(4)	89.49(3)	92.13(1)	94.24(10)	98.41(5)	96.01(8)	95.25(9)	84.32(15)	97.37(3)	99.22(3)	97.91(6)	98.93(4)	99.22(2)	99.62(1)		
Hepatitis	58.17(16)	95.32(14)	98.57(6)	98.11(7)	97.06(9)	96.81(11)	98.75(4)	96.54(12)	64.77(15)	98.79(3)	96.16(13)	97.03(10)	98.64(5)	99.22(2)	99.48(1)			
http	99.80(11)	99.99(7)	99.98(8)	100.00(1)	100.00(1)	99.96(10)	98.33(12)	98.33(12)	83.31(14)	0.30(16)	99.00(1)	100.00(5)	43.69(15)	100.00(1)	99.79(9)			
InternetAds	68.43(14)	81.07(9)	90.28(2)	84.65(7)	82.06(8)	80.87(10)	86.62(6)	65.24(15)	80.36(11)	79.98(12)	74.64(13)	N/A(NA)	87.18(3)	86.81(5)	86.87(4)	92.63(1)		
lensesphere	55.52(15)	91.47(9)	93.75(11)	98.76(4)	99.59(8)	99.81(6)	98.75(1)	89.03(13)	89.94(12)	90.50(11)	97.50(14)	98.32(6)	98.43(2)	98.39(1)	97.67(3)			
landesk	50.80(16)	92.87(3)	64.29(15)	71.09(12)	68.76(12)	82.95(11)	92.54(11)	71.89(11)	83.09(13)	90.59(8)	85.30(9)	92.26(6)	91.88(6)	92.25(4)	91.73(7)	94.47(1)		
letter	69.93(25)	78.65(10)	82.72(6)	78.91(9)	81.17(7)	73.01(14)	90.33(2)	61.36(16)	73.85(12)	75.00(12)	86.46(3)	77.73(1)	85.75(4)	85.20(5)	80.78(8)	96.38(1)		
Lymphography	97.21(11)	97.27(10)	99.80(2)	99.93(6)	97.24(9)	94.91(12)	90.21(15)	99.85(1)	77.53(16)	94.34(7)	99.88(13)	99.58(5)	99.14(8)	98.61(9)	99.79(3)			
magic_gamma	50.32(16)	88.02(5)	81.14(4)	82.70(11)	82.91(10)	84.14(9)	89.24(2)	78.10(14)	65.59(15)	89.45(1)	81.61(12)	86.72(7)	88.40(4)	87.64(6)	88.85(3)			
mammography	75.96(15)	94.46(1)	92.81(5)	92.78(6)	94.14(3)	92.47(10)	99.12(2)	95.07(15)	98.15(19)	95.00(13)	97.13(1)	89.85(14)	97.78(6)	98.74(3)	98.79(1)			
mnist	68.76(16)	96.65(12)	98.89(4)	98.77(5)	98.93(4)	98.03(13)	98.96(11)	99.14(1)	98.42(12)	99.06(10)	98.71(13)	98.19(1)	98.67(7)	98.46(8)	99.20(1)			
mus	90.41(16)	100.00(10)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)		
optdigits	46.49(16)	99.87(8)	99.92(5)	99.99(3)	99.97(4)	99.96(3)	99.95(2)	97.64(14)	96.62(14)	99.94(5)	90.00(15)	99.63(6)	99.72(10)	99.62(12)	99.73(11)	99.88(7)		
PageBlocks	64.80(16)	96.18(5)	94.42(10)	88.59(4)	90.15(3)	95.31(7)	97.65(3)	91.16(12)	85.81(15)	93.45(1)	94.69(9)	96.09(6)	95.05(8)	97.60(4)	97.72(2)	98.36(1)		
pima	61.46(16)	73.27(12)	71.75(13)	81.52(4)	81.06(5)	76.66(11)	82.55(3)	80.49(8)	63.21(15)	76.15(15)	77.00(10)	86.15(5)	77.36(9)	80.87(6)	83.38(1)	80.78(7)	82.56(2)	
satellite	74.94(16)	93.57(5)	79.72(4)	84.67(1)	84.71(10)	85.89(9)	93.21(3)	78.73(15)	80.17(13)	90.88(8)	84.07(10)	93.64(4)	93.31(6)	92.74(7)	94.37(1)			
satisfaction-2	90.63(16)	92.26(1)	91.68(2)	92.86(1)	92.91(3)	91.96(1)	92.92(1)	92.92(1)	92.93(1)	92.93(1)	92.93(1)	92.93(1)	92.93(1)	92.93(1)	92.93(1)	92.93(1)		
shuttle	79.04(16)	89.85(6)	98.94(3)	97.57(2)	97.67(4)	98.67(3)	99.90(1)	99.90(1)	79.26(14)	88.54(5)	97.57(13)	98.60(9)	98.35(6)	98.15(9)	99.77(2)	99.54(3)		
skin	52.75(16)	99.88(2)	89.70(5)	95.99(5)	92.57(3)	98.02(10)	98.94(4)	93.96(14)	88.96(12)	92.32(11)	97.90(6)	99.84(3)	99.69(4)	99.67(9)	99.69(7)	99.91(1)		
smtp	51.69(14)	82.74(3)	76.47(8)	75.09(9)	75.24(1)	75.41(10)	79.87(5)	62.05(10)	98.60(12)	86.40(1)	72.61(12)	97.28(16)	95.51(7)	96.51(7)	94.32(11)	96.70(6)	96.40(8)	
SpamBase	55.09(16)	90.60(10)	89.77(1)	91.69(7)	91.03(8)	92.05(11)	97.02(11)	93.55(19)	92.15(14)	92.97(17)	96.97(16)	99.77(6)	99.72(16)	99.72(14)	99.70(16)	99.70(14)	99.70(14)	
WBC	75.93(13)	91.04(11)	91.03(10)	91.04(1)	91.03(11)	91.03(10)	91.04(11)	91.04(11)	91.04(11)	91.04(11)	91.04(11)	91.04(11)	91.04(11)	91.04(11)	91.04(11)	91.04(11)		
Wilt	36.73(16)	96.83(8)	92.64(15)	87.97(1)	87.83(6)	85.66(10)	84.53(4)	87.83(1)	80.50(10)	82.25(10)	85.02(15)	80.65(13)	83.40(12)	83.33(14)	90.97(16)	99.72(16)	96.30(1)	
wine	66.97(16)	97.51(16)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	
WPBC	47.93(16)	80.12(10)	73.53(4)	79.68(1)	80.45(8)	78.50(2)	87.98(3)	67.98(5)	70.13(13)	80.32(9)	90.04(1)	84.27(7)	85.86(5)	87.06(5)	87.80(5)	87.28(5)	94.41(1)	
yeast	48.40(15)	71.47(7)	66.29(8)	66.76(4)	67.70(3)	69.64(8)	64.52(12)	98.07(10)	59.53(14)	59.53(15)	60.51(13)	62.36(13)	61.30(13)	66.70(5)	65.98(9)	72.44(1)		
CIFAR10	60.91(14)	70.82(11)	30.87(6)	82.32(2)	76.80(9)	77.06(7)	81.58(3)	53.08(16)	76.43(9)	73.54(1)	59.59(5)	64.30(2)	62.01(13)	80.67(5)	81.27(4)	83.26(1)		
FashionMNIST	80.05(14)	79.31(9)	49.52(4)	95.72(1)	94.21(7)	94.21(7)	93.89(8)	95.23(6)	61.06(1)	94.20(12)	73.16(15)	73.16(15)	98.19(10)	98.19(10)	98.19(10)	98.19(10)	94.29(6)	
MNISTC	7.84(16)	91.21(10)	92.11(1)	92.21(1)	93.45(1)	94.06(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	95.23(1)	
MTTC-AD	47.45(14)	88.14(1)	45.47(1)	89.73(5)	91.26(5)	52.16(3)	52.16(3)	19.42(15)	25.74(14)	24.39(11)	45.08(8)	57.69(2)	60.68(1)	0.14(16)	48.75(7)	50.16(6)		
MTTC-AD	48.74(16)	87.78(15)	45.77(1)	88.79(1)	89.20(8)	89.31(9)	90.34(1)	51.35(2)	49.85(12)	99.00(3)	84.30(9)	92.14(12)	92.28(12)	94.22(12)	94.42(12)	95.17(1)		
SVHN	7.93(15)	16.07(12)	24.51(3)	27.06(1)	26.68(2)	19.72(9)	5.19(16)	20.64(6)	18.74(10)	19.72(9)	9.83(14)	20.97(15)	19.53(4)	88.16(10)	88.77(9)	90.25(1)		
Agnews	7.28(16)	52.97(6)	53.35(5)	54.60(4)	55.92(3)	57.65(2)	45.38(9)	8.56(15)	46.99(8)	64.87(1)	42.92(11)	38.17(13)	11.38(4)	10.46(12)	10.46(12)	11.38(4)	85.03(1)	
Amazon	6.04(16)	16.95(12)	27.75(4)	30.39(2)	29.65(3)	25.77(5)	22.95(7)	5.41(16)	12.42(13)	23.11(14)	11.38(14)	33.09(4)	6.00(14)	19.19(10)	23.89(6)			
Imdb	4.83(16)	18.16(12)	21.82(5)	23.03(3)	23.47(5)	27.31(2)	22.66(4)	5.82(16)	23.46(14)	33.01(1)	23.14(1)	13.84(13)	6.06(14)	20.68(11)	11.27(7)	23.07(7)		
Yelp	9.00(14)	23.29(13)	19.42(2)	19.80(1)	19.09(4)	23.59(2)	35.67(8)	5.81(15)	21.22(10)	31.82(9)	33.08(1)	20.04(9)	8.84(15)	30.88(10)	33.75(8)	34.80(7)		
20news	7.28(14)	18.37(11)	7.20(15)	28.30(3)	30.92(2)	26.57(4)	22.26(2)	5.26(16)	20.84(9)	34.04(1)	21.36(8)	14.13(13)	15.73(12)	21.60(7)</td				

Table D14: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 50\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRENet	FEAWAD	XGBOD	NB	SVM	MLP	ResNet	FFTransformer	RF	LGB	XGB	CatB	
ALOI	54.52(1)	64.14(6)	53.63(12)	52.88(14)	52.95(8)	57.97(5)	79.27(1)	52.70(15)	52.07(16)	57.46(8)	56.16(9)	54.72(10)	78.90(2)	70.12(5)	70.28(4)		
anthyroid	76.85(16)	95.54(9)	83.31(11)	82.81(3)	82.58(4)	91.43(10)	99.21(5)	83.01(12)	82.03(15)	96.62(8)	98.13(7)	98.87(6)	99.34(3)	99.38(1)	99.37(2)	99.27(4)	
backdoor	86.48(15)	98.24(6)	89.60(14)	98.11(7)	94.07(12)	98.48(5)	97.27(9)	85.15(6)	97.21(10)	97.69(8)	92.79(13)	96.56(11)	99.69(1)	99.23(4)	99.31(3)	99.64(2)	
breastw	94.15(16)	98.44(13)	98.84(12)	99.67(2)	99.59(4)	98.97(1)	99.17(9)	99.75(1)	98.32(14)	99.62(3)	95.52(15)	99.10(10)	99.19(8)	99.42(5)	99.40(6)	99.38(7)	
campaign	58.36(15)	78.11(13)	57.94(16)	86.53(8)	87.42(7)	77.21(4)	92.93(1)	81.14(11)	78.45(12)	82.95(10)	84.07(9)	89.47(6)	91.25(2)	92.46(4)	92.17(1)	92.86(2)	
carbs	85.53(16)	96.23(6)	96.38(7)	98.03(6)	98.07(5)	96.56(3)	99.93(1)	94.03(10)	95.22(10)	96.80(12)	96.66(11)	98.27(1)	99.29(2)	99.31(3)	99.43(5)		
Cardiotocography	62.27(6)	94.69(11)	95.80(10)	95.05(10)	95.22(0)	94.02(8)	96.00(4)	93.07(14)	92.20(9)	95.79(7)	98.09(15)	98.55(3)	96.23(2)	96.42(2)	95.95(5)	95.95(1)	
celiba	69.88(15)	92.06(11)	57.64(16)	95.71(3)	95.34(6)	93.11(9)	95.81(2)	89.48(11)	95.44(4)	82.83(4)	95.42(5)	89.42(12)	95.91(1)	93.73(8)	95.04(7)		
census	58.46(16)	82.01(12)	69.47(17)	91.36(4)	89.89(7)	84.60(11)	92.51(2)	73.31(13)	86.63(9)	70.73(4)	86.68(8)	90.87(6)	91.72(3)	91.29(5)	92.95(1)		
cover	42.11(16)	99.95(3)	99.95(1)	99.93(7)	99.94(6)	99.88(8)	99.64(13)	98.91(12)	98.99(14)	99.86(6)	97.07(5)	99.94(4)	99.95(2)	99.83(10)	99.82(11)	99.94(5)	
donors	57.08(16)	100.00(1)	82.81(10)	100.00(1)	99.94(4)	99.99(5)	99.98(4)	99.60(14)	99.91(11)	100.00(1)	99.77(13)	100.00(1)	99.88(12)	99.98(7)	99.96(9)	99.98(6)	
fault	66.23(16)	79.07(5)	76.11(15)	77.89(2)	77.47(9)	75.53(13)	81.29(4)	69.99(15)	87.19(12)	79.71(4)	78.04(8)	83.16(11)	88.90(1)	90.15(6)	84.04(1)		
glass	63.48(16)	95.63(10)	90.11(11)	89.00(4)	90.35(12)	99.61(6)	99.66(5)	94.38(11)	89.52(9)	85.43(12)	99.83(2)	99.67(4)	99.73(3)	99.34(8)	99.41(7)	99.87(1)	
Hepatitis	66.13(16)	99.67(11)	99.74(10)	99.23(12)	98.83(3)	99.96(8)	97.18(4)	94.81(15)	100.00(1)								
http	99.77(11)	100.00(1)	99.98(8)	100.00(1)	100.00(1)	99.93(10)	98.32(12)	98.33(14)	0.39(16)	98.49(1)	98.53(10)	80.04(13)	N/A(N/A)	93.57(3)	92.93(6)	93.48(4)	
InternetAds	68.93(16)	90.13(9)	94.33(2)	92.95(5)	92.47(7)	89.31(2)	92.35(8)	76.24(4)	89.42(1)	90.98(1)	94.26(6)	89.24(11)	94.26(10)	94.26(6)	94.26(5)	94.26(1)	
ionsphere	92.90(16)	96.53(15)	97.91(4)	92.22(6)	92.36(1)	98.41(5)	98.73(6)	98.74(9)	98.20(16)	98.37(7)	99.05(8)	98.77(5)	98.77(2)	98.77(3)	98.77(4)	98.77(1)	
landsat	53.50(16)	94.52(4)	60.57(15)	90.19(2)	90.66(1)	90.64(11)	94.60(8)	71.19(14)	92.98(8)	90.21(12)	90.39(7)	90.39(5)	92.89(8)	91.23(2)	90.07(4)	90.45(5)	
letter	70.16(16)	83.75(10)	88.74(6)	84.17(9)	83.32(11)	91.91(2)	70.70(5)	81.80(4)	82.59(2)	84.90(8)	91.35(5)	90.57(4)	91.58(3)	87.40(7)	89.45(2)	94.25(1)	
Lymphography	98.18(14)	99.88(13)	100.00(1)	100.00(1)	100.00(1)	99.95(4)	99.96(5)	99.97(6)	99.98(7)	99.98(8)	99.99(9)	99.99(10)	99.99(11)	99.99(12)	99.99(13)	99.99(14)	
magic_gamma	56.32(19)	90.24(4)	84.26(14)	83.19(12)	84.89(10)	90.13(6)	77.43(5)	83.81(11)	90.41(3)	85.29(9)	90.68(2)	90.18(2)	90.48(3)	91.24(7)	90.18(5)	91.07(1)	
mammography	77.28(15)	94.42(2)	93.17(14)	93.28(8)	93.29(8)	93.90(5)	91.32(2)	70.70(6)	93.28(9)	94.12(6)	94.12(7)	93.92(8)	94.12(8)	94.48(1)			
mnist	67.77(16)	99.14(7)	99.36(7)	99.01(10)	99.13(8)	97.21(3)	90.02(5)	98.01(12)	98.42(11)	95.66(4)	99.10(9)	99.56(4)	99.26(4)	99.31(5)	99.19(6)	99.51(2)	
mask	90.89(16)	99.37(10)	99.59(9)	99.46(12)	99.51(8)	99.85(14)	99.96(13)	99.97(15)	99.98(16)	99.99(17)	99.99(18)	99.99(19)	99.99(20)	99.99(21)	99.99(22)	99.99(23)	
spiders	46.90(16)	99.97(6)	99.99(3)	99.99(4)	99.99(5)	99.99(6)	99.99(7)	98.05(14)	98.89(13)	99.98(5)	99.98(6)	99.91(11)	99.91(10)	99.92(10)	99.92(9)	99.92(8)	
PageBlocks	78.46(16)	97.21(8)	94.08(12)	88.09(5)	90.53(4)	96.22(9)	98.31(4)	91.83(13)	94.17(11)	95.63(10)	97.49(7)	97.61(6)	98.29(5)	98.39(2)	98.34(3)	98.87(1)	
pendigits	59.36(16)	99.99(2)	99.88(8)	99.98(9)	99.99(11)	99.99(12)	99.99(13)	99.99(14)	99.99(15)	99.99(16)	99.99(17)	99.99(18)	99.99(19)	99.99(20)	99.99(21)	100.00(1)	
Pima	63.36(16)	80.63(11)	75.75(15)	82.93(9)	82.73(10)	80.50(12)	88.62(5)	83.54(7)	77.38(4)	82.97(8)	80.45(5)	84.25(6)	89.73(1)	89.63(2)	88.73(4)	89.13(3)	
satellite	75.12(16)	95.28(7)	81.26(1)	85.38(5)	83.62(1)	95.31(5)	98.87(5)	78.88(15)	91.43(8)	90.66(10)	91.39(7)	95.39(6)	95.25(1)	95.42(6)	95.40(5)	96.19(2)	
satimage-2	96.83(16)	98.21(7)	98.94(6)	98.77(12)	98.61(1)	98.84(7)	98.87(7)	98.94(11)	98.97(12)	98.97(13)	98.97(14)	98.97(15)	98.97(16)	98.97(17)	98.97(18)	98.97(19)	
shuttle	65.45(16)	98.97(6)	98.91(7)	97.57(14)	97.67(1)	98.47(2)	98.31(2)	97.37(15)	98.28(10)	97.63(13)	98.68(9)	98.28(11)	99.52(5)	99.65(4)	99.92(2)		
skin	52.08(16)	99.92(3)	98.03(2)	95.70(2)	95.23(8)	98.77(11)	99.88(7)	93.94(14)	99.73(14)	99.88(6)	99.93(2)	99.90(4)	99.77(10)	99.89(5)	99.87(8)	99.94(1)	
smtip	56.11(14)	99.22(4)	92.43(6)	87.46(7)	84.95(9)	94.22(11)	94.23(12)	94.24(13)	95.04(11)	95.05(12)	95.36(13)	95.36(14)	95.36(15)	95.36(16)	95.36(17)		
SpamBase	57.07(16)	94.37(8)	91.77(1)	92.05(2)	94.09(7)	94.93(1)	90.40(12)	88.06(12)	88.30(13)	90.45(12)	90.75(13)	90.75(14)	90.75(15)	90.75(16)	90.75(17)		
speech	47.72(16)	95.94(14)	73.34(7)	80.36(2)	81.56(1)	75.98(6)	74.21(2)	77.34(8)	79.76(5)	72.88(8)	79.25(4)	86.81(5)	86.69(12)	86.35(13)	86.34(13)		
spiders	93.97(16)	98.63(1)	99.05(2)	97.96(3)	99.35(4)	97.97(5)	99.27(6)	97.97(7)	99.27(8)	97.97(9)	99.27(10)	97.97(11)	99.27(12)	99.27(13)	99.27(14)		
thyroid	94.52(16)	99.04(13)	97.76(5)	97.96(6)	99.90(1)	97.97(8)	97.98(9)	97.99(10)	97.99(11)	97.99(12)	97.99(13)	97.99(14)	97.99(15)	97.99(16)	97.99(17)		
vertebral	38.89(16)	88.40(10)	79.82(1)	81.09(3)	82.54(2)	92.65(8)	98.71(1)	87.03(11)	91.13(9)	74.28(15)	96.42(6)	97.39(7)	97.95(8)	97.87(4)	98.22(3)	98.59(2)	
wovels	78.97(16)	93.21(14)	98.94(7)	98.08(9)	98.94(3)	98.13(7)	98.13(8)	97.53(11)	98.14(6)	97.40(13)	98.20(15)	99.60(1)	96.27(12)	98.85(5)	97.60(10)	99.43(2)	
Waveform	53.02(16)	85.40(4)	91.74(4)	46.79(6)	47.51(7)	95.37(3)	91.05(6)	86.98(12)	91.97(2)	94.73(4)	75.22(15)	90.30(9)	87.11(11)	90.44(8)	86.91(10)	91.77(3)	
WBC	95.21(16)	98.38(10)	97.81(1)	95.41(4)	95.92(5)	99.44(1)	92.44(1)	98.94(1)	99.45(2)	99.60(1)	95.20(1)	97.31(1)	97.46(1)	98.95(3)	95.94(1)		
WDBC	97.77(16)	99.72(12)	99.88(1)	99.92(1)	99.93(1)	99.94(1)	99.95(1)	99.96(1)	99.97(1)	99.98(1)	99.99(1)	99.99(1)	99.99(1)	99.99(1)	99.99(1)		
Wilt	34.40(16)	96.93(5)	96.35(6)	96.18(14)	96.10(12)	96.23(11)	96.12(1)	122(24)	37.34(8)	36.76(9)	27.78(8)	33.71(11)	48.97(4)	49.10(3)	47.54(5)	53.74(2)	
wine	70.50(16)	99.89(14)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	98.73(15)	87.81(10)	91.62(9)	83.88(13)	92.57(7)	97.57(5)	80.36(15)	89.69(11)	91.82(8)		
WPBC	49.03(16)	89.94(8)	94.10(2)	78.26(10)	81.43(12)	95.84(4)	82.55(1)	80.90(14)	98.51(13)	98.19(14)	98.18(15)	99.20(16)	99.35(17)	99.73(18)			
yeast	48.87(16)	74.61(10)	64.07(11)	84.17(14)	85.42(12)	85.45(14)	87.86(7)	69.93(15)	86.82(6)	75.90(11)	80.67(8)	85.77(3)	84.85(6)	81.78(7)	86.87(1)		
ZF	74.88(16)	97.82(12)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	98.71(15)	94.16(13)	94.16(14)	94.16(15)	94.16(16)	94.16(17)	94.16(18)	94.16(19)	94.16(20)		
CIFAR10	92.7(15)	97.2(1)	97.92(1)	35.42(2)	34.62(5)	28.25(9)	35.17(4)	74.41(6)	28.54(8)	26.97(10)	23.86(2)	19.38(3)	14.98(14)	34.03(7)	34.29(6)	35.23(3)	
EmotionMNIST	21.70(16)	73.78(18)	81.48(7)	36.52(3)	80.54(6)	71.13(10)	72.22(1)	22.78(1)	67.32(12)	61.52(12)	59.75(13)	41.20(13)	80.39(14)	91.92(11)	85.77(5)		
MNIST-C	19.26(16)	32.75(10)	42.14(14)	86.63(1)	86.52(10)	82.10(10)	83.19(8)	26.52(15)	53.53(16)	44.64(15)	92.59(7)	91.02(12)	90.40(13)	91.35(15)			
MVtec-AD	58.20(14)	86.43(6)	55.61(15)	83.77(8)	81.43(12)	22.08(11)	40.77(4)	14.57(15)	57.61(1)	42.15(3)	36.74(8)	36.86(7)	35.78(9)	37.47(6)	37.65(5)	42.67(2)	
SVHN	7.95(15)	22.28(12)	30.55(3)</td														

Table D16: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 75\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Datasets	GANomaly	DeepSAD	REPEN	DevNet	PRNet	FEAWD	XGBOD	NB	SVM	MLP	ResNet	FTTransformer	RF	LGB	XGB	CatB
ALOI	56.22(11)	66.49(6)	51.94(15)	52.85(14)	50.82(16)	63.28(7)	86.73(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)	98.10(8)	97.96(9)	98.61(7)	99.10(6)	99.46(2)	99.44(3)	99.47(1)	99.42(4)
backdoor	86.89(16)	98.69(6)	89.62(15)	97.91(9)	94.93(13)	98.56(8)	97.52(10)	93.90(14)	98.56(7)	98.88(5)	96.33(12)	96.96(11)	99.95(1)	99.73(2)	99.60(4)	99.64(3)
breastw	95.04(16)	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)	99.70(8)	99.79(1)	99.70(7)	99.71(6)	99.78(2)
campaign	66.24(15)	81.12(13)	58.00(16)	87.18(9)	88.34(7)	76.01(14)	93.54(1)	81.89(12)	84.32(11)	86.75(10)	88.08(3)	90.05(6)	93.84(5)	93.29(3)	92.47(4)	93.52(2)
cardio	87.03(16)	96.71(11)	90.86(10)	98.09(7)	94.14(12)	74.84(11)	99.74(1)	94.77(10)	96.73(10)	96.73(10)	96.73(10)	96.73(10)	99.42(4)	99.91(3)	99.92(5)	99.92(3)
Cardiotocography	66.60(16)	96.59(8)	95.10(10)	85.15(11)	95.25(2)	96.04(1)	97.60(4)	94.74(14)	96.57(9)	96.92(7)	94.84(11)	96.45(4)	97.96(2)	97.79(3)	97.45(6)	97.45(1)
celeba	63.72(15)	93.78(10)	57.64(15)	95.82(3)	95.78(4)	93.98(9)	96.28(1)	89.66(13)	96.57(9)	95.47(6)	86.43(4)	95.43(7)	90.59(12)	96.18(2)	95.15(8)	95.74(5)
census	63.95(16)	85.99(12)	69.68(15)	92.33(5)	91.29(7)	87.55(10)	73.51(4)	88.61(9)	89.06(8)	79.60(3)	87.17(11)	92.22(6)	92.85(3)	92.58(4)	93.55(1)	92.58(4)
cover	42.57(16)	99.98(1)	99.95(7)	99.98(6)	99.80(13)	99.72(15)	99.89(12)	99.97(2)	99.97(5)	99.95(6)	99.94(11)	99.96(4)	99.95(10)	99.78(14)	99.97(3)	99.97(3)
donors	56.37(16)	100.00(1)	82.83(15)	100.00(1)	99.95(5)	99.96(2)	100.00(1)	99.60(14)	100.00(1)	100.00(1)	99.83(13)	100.00(1)	99.99(10)	100.00(7)	99.99(11)	100.00(8)
fault	67.43(16)	81.85(12)	77.92(13)	79.84(14)	75.71(12)	86.71(15)	83.85(7)	69.77(15)	80.31(7)	78.37(10)	82.06(7)	86.30(7)	83.86(7)	86.93(1)	86.93(1)	86.93(1)
fraud	9.10(16)	92.66(6)	90.93(9)	22.38(7)	32.62(5)	9.34(1)	96.68(1)	89.87(11)	87.08(15)	85.67(13)	87.08(10)	84.45(4)	88.05(12)	87.81(2)	90.55(1)	90.55(1)
glass	68.70(16)	97.73(10)	88.86(15)	88.60(15)	90.29(9)	95.68(12)	93.98(8)	93.74(11)	99.03(12)	98.50(12)	99.52(12)	99.99(7)	99.99(6)	99.99(6)	99.99(6)	99.99(6)
Hepatitis	69.27(16)	99.92(9)	99.78(11)	98.86(14)	99.36(2)	99.14(13)	99.93(8)	97.10(15)	99.85(10)	100.00(1)						
http	99.80(12)	100.00(1)	99.98(9)	100.00(1)												
InternetAids	69.50(15)	93.90(11)	96.22(2)	95.92(4)	95.33(14)	94.86(7)	91.40(6)	95.69(15)	95.69(15)	96.68(13)	94.09(14)	98.08(3)	NA(NA)	94.89(6)	94.47(8)	95.14(5)
ionsphere	7.25(16)	7.32(15)	6.25(14)	6.89(13)	6.94(12)	6.84(11)	7.44(10)	7.34(15)	7.47(13)	7.46(14)	7.49(12)	7.51(13)	7.51(12)	7.51(13)	7.51(13)	7.51(13)
landuse	5.72(16)	9.27(15)	6.55(14)	6.46(13)	6.85(14)	6.25(11)	5.98(10)	6.22(15)	6.23(14)	6.19(13)	6.21(14)	6.22(13)	6.22(13)	6.22(13)	6.22(13)	6.22(13)
letter	70.42(16)	86.83(11)	88.48(8)	87.56(7)	86.44(6)	84.29(4)	93.35(5)	88.48(9)	87.60(9)	91.66(6)	90.85(7)	95.39(2)	93.39(4)	93.86(3)	96.49(1)	93.86(3)
Lymphography	98.31(16)	100.00(1)														
magic_gamma	85.56(16)	83.16(12)	83.14(14)	83.19(13)	84.17(11)	90.97(4)	77.08(15)	90.47(8)	88.30(13)	89.89(12)	91.88(11)	91.38(12)	91.88(11)	91.38(12)	92.25(1)	92.25(1)
mammography	75.53(16)	95.71(12)	93.01(9)	93.08(7)	92.81(12)	95.62(10)	92.81(12)	81.77(15)	93.34(7)	92.92(15)	92.52(15)	93.46(6)	90.41(14)	92.67(13)	95.67(3)	95.67(3)
mnist	68.69(16)	99.93(10)	98.99(11)	98.99(11)	98.99(12)	98.76(12)	99.87(1)	94.97(15)	99.32(13)	99.34(12)	99.34(11)	99.34(11)	99.72(6)	99.55(5)	99.71(3)	99.71(3)
mask	100.00(16)	100.00(1)														
optdigits	47.53(16)	99.99(2)	99.98(1)	99.99(1)	99.98(10)	99.98(8)	99.98(1)	98.50(9)	98.51(10)	98.52(9)	98.53(10)	98.53(10)	98.53(10)	98.53(10)	98.53(10)	98.53(10)
PageBlocks	79.45(16)	98.08(7)	95.25(12)	95.10(14)	95.10(14)	95.55(11)	95.55(11)	94.10(1)	95.55(11)	95.55(11)	95.55(11)	95.55(11)	95.55(11)	95.55(11)	95.55(11)	95.55(11)
pendigits	62.56(16)	99.99(2)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)	99.98(1)
Pima	65.20(16)	83.81(10)	75.97(15)	82.15(25)	83.38(12)	82.79(9)	90.47(5)	83.47(11)	86.61(7)	84.51(9)	85.94(8)	88.07(6)	93.84(1)	92.00(3)	91.30(4)	92.02(2)
satellite	77.11(16)	96.22(5)	81.55(12)	84.27(13)	84.99(12)	86.32(13)	92.03(7)	78.66(15)	94.43(9)	94.43(9)	94.66(6)	95.88(7)	97.11(1)	96.76(9)	96.46(4)	96.98(2)
satimage-2	97.08(16)	96.93(14)	96.92(14)	96.93(13)	96.92(14)	96.93(13)	96.93(13)	96.93(13)	96.93(13)	96.93(13)	96.93(13)	96.93(13)	97.01(16)	97.01(16)	97.01(16)	97.01(16)
shuttle	78.80(16)	99.43(6)	98.69(9)	97.57(14)	97.59(11)	98.21(11)	99.99(5)	97.74(12)	88.89(7)	97.56(15)	97.56(15)	100.00(2)	100.00(1)	99.99(3)	99.99(4)	99.99(3)
skin	52.94(16)	99.92(8)	89.30(12)	95.77(12)	95.28(12)	93.77(11)	99.71(1)	93.90(14)	99.94(5)	99.96(3)	99.99(3)	99.99(3)	99.99(3)	99.99(3)	99.99(3)	99.99(3)
smtp	55.30(16)	99.23(5)	91.98(8)	85.48(11)	92.44(9)	79.88(15)	92.45(3)	80.38(14)	98.81(14)	90.46(9)	90.46(10)	90.46(10)	95.76(1)	92.15(1)	92.15(1)	92.15(1)
SpamBase	59.25(16)	95.94(10)	94.20(12)	93.02(13)	94.01(12)	93.02(13)	95.81(1)	11.20(14)	20.95(11)	20.89(9)	23.20(13)	30.74(7)	29.12(10)	34.23(2)	30.63(8)	34.07(4)
speech	47.76(16)	61.41(4)	79.23(6)	84.01(2)	83.52(1)	75.16(10)	77.09(9)	74.49(11)	80.23(8)	81.29(4)	73.85(12)	77.27(8)	60.60(15)	69.39(13)	79.68(5)	77.38(7)
task	93.27(16)	99.55(15)	99.80(9)	99.87(11)	99.98(12)	99.98(13)	99.99(1)	99.65(13)	99.77(10)	99.72(12)	99.59(11)	99.56(10)	99.94(13)	99.94(13)	99.94(13)	99.94(13)
thyroid	44.27(16)	86.81(12)	86.70(12)	86.80(12)	86.80(12)	86.80(12)	86.86(1)	88.73(11)	92.19(12)	92.19(12)	92.19(12)	92.19(12)	92.19(12)	92.19(12)	92.19(12)	92.19(12)
vertebral	40.75(16)	82.04(13)	66.75(9)	66.75(13)	66.75(13)	66.75(13)	67.02(1)	55.18(13)	55.18(13)	55.18(13)	55.18(13)	55.18(13)	55.18(13)	55.18(13)	55.18(13)	55.18(13)
vowels	40.75(16)	81.76(14)	81.74(15)	82.32(13)	93.53(8)	99.14(3)	98.43(1)	84.63(12)	92.39(7)	98.87(9)	97.56(15)	100.00(2)	98.88(1)	99.93(2)	99.93(2)	99.93(2)
Waveform	44.74(16)	47.70(6)	97.72(15)	98.83(14)	98.83(14)	98.83(12)	99.67(1)	99.37(10)	98.87(9)	98.76(13)	98.57(12)	98.41(12)	98.47(10)	98.76(11)	98.76(11)	98.76(11)
WBC	48.77(16)	88.02(11)	89.93(11)	94.77(9)	84.18(15)	95.14(7)	95.95(1)	86.97(11)	94.59(14)	94.59(14)	94.59(14)	94.59(14)	94.59(14)	94.59(14)	94.59(14)	94.59(14)
WDBC	61.00(16)	99.65(11)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)	100.00(1)
Wile	4.81(16)	84.42(9)	7.11(15)	8.19(14)	8.49(13)	55.31(10)	88.49(3)	40.55(11)	87.01(11)	48.63(12)	24.46(12)	87.91(6)	98.92(11)	88.82(8)	88.82(1)	88.37(7)
wine	24.84(16)	99.85(15)	100.00(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(4)
yeast	33.54(16)	51.46(11)	36.34(15)	49.24(12)	48.99(13)	54.83(7)	60.76(5)	46.99(14)	52.52(9)	51.54(10)	55.07(6)	54.38(8)	63.83(2)	61.78(3)	61.63(4)	65.46(1)
CIFAR10	94.2(16)	31.80(11)	41.15(1)	38.42(4)	36.93(8)	33.51(9)	38.04(7)	9.52(15)	38.16(6)	32.83(10)	30.76(12)	24.34(13)	19.97(14)	39.16(3)	38.24(5)	39.37(2)
EpsilonMNIST	22.73(16)	81.89(9)	82.56(6)	82.73(6)	82.73(6)	83.00(11)	84.95(10)	26.24(10)	88.23(8)	78.71(1)	70.83(1)	72.97(1)	84.12(12)	84.12(12)	84.12(12)	84.12

Table D18: AUCROC of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 100\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.

Table D19: AUCPR of 16 label-informed algorithms on 57 benchmark datasets, with labeled anomaly ratio $\gamma_l = 100\%$. We show the performance rank in parenthesis (lower the better), and mark the best performing method(s) in **bold**.