

# Comparison of Anomaly Detectors: Context Matters

IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS

Special Issue on Deep Learning for Anomaly Detection

Vít Škvára, Jan Franců, Matěj Zorek, Tomáš Pevný, *Member, IEEE*, Václav Šmídl, *Member, IEEE*

**Abstract**—Deep generative models are challenging the classical methods in the field of anomaly detection nowadays. Every new method provides evidence of outperforming its predecessors, often with contradictory results. The objective of this comparison is twofold: to compare anomaly detection methods of various paradigms with focus on deep generative models, and identification of sources of variability that can yield different results. The methods were compared on popular tabular and image datasets. We identified the main sources of variability to be experimental conditions: i) the type data set (tabular or image) and the nature of anomalies (statistical or semantic), and ii) strategy of selection of hyperparameters, especially the number of available anomalies in the validation set. Different methods perform the best in different contexts, i.e. combination of experimental conditions together with computational time. This explains the variability of the previous results and highlights the importance of careful specification of the context in the publication of a new method. All our code and results are available for download.

## I. INTRODUCTION

DEEP generative models are gaining popularity in anomaly detection since the introduction of the Variational Autoencoder (VAE) [45]. The number of modifications and extensions of VAE or generative adversarial networks (GAN) [33] is sharply increasing, each claiming superiority over the prior art. This raises a suspicion that some of the methods are overspecialized or poorly tested. This work, inspired by the paper “Do we need hundreds of classifiers to solve real-world classification problems?” [29], strives to compare under “fair” conditions anomaly detectors to observe how the field has evolved in the last twenty years (the oldest compared detector [73] was published in 2000). Specifically, it investigates if methods based on *deep* generative models offer a benefit over methods based on alternative paradigms, either the *classical* methods based on distances, or deep architectures without the capability of generating samples.

Surely there already exist comparisons of anomaly detectors. Earlier surveys [71, 12, 32, 68] do not compare to deep generative methods, because they were not developed or sufficiently popular at that time. Contrary to that, the study in [46] contains a detailed description of deep models, but provides experiments only with the basic VAE and only on specialized video datasets. Ref. [13] introduces a taxonomy of deep anomaly detection models but does not compare them experimentally. Other recent surveys [57, 52, 28, 88, 62] either ignore deep generative models altogether or describe

them only theoretically, without making any experimental comparison. The most relevant prior art is [75], which tries to theoretically link deep and shallow techniques<sup>1</sup>. But again, an extensive experimental comparison of different generative models is missing. One would also expect papers introducing new methods to contain such a comparison. Some of them do [68], but generally, we have found comparisons limited (e.g. using a small number of datasets or methods) or flawed, which is elaborated below.

How does this paper avoid the aforementioned deficiencies? First, eight classical methods in comparison serve as a baseline, over which we expect the state-of-the-art deep methods should improve upon (latest compared method [89] was published in 2020). Second, the comparison uses a large number of tabular (40) and image (6) datasets popular in the evaluation of deep models. Third, all methods have been given the same conditions, which primarily means optimization of hyperparameters, as [82] has shown this to have a significant impact.

The list of contributions contains:

- 1) Experimental comparison of classical and deep anomaly detectors on a large number of datasets.
- 2) Identification of the dataset type, hyper-parameter selection strategy, and computational cost as major factors in the selection of the most suitable method.
- 3) We publish codes of the evaluation pipeline and compared methods, including automatic download of datasets, splitting them into training, validation, and testing, and calculating the performance metrics.

The paper is organized as follows. First, the anomaly detection contexts that had the greatest influence on the outcome of our experiments are defined. Second, there is a brief theoretical overview of the tested generative deep models and other methods. Sec. IV details the datasets, different approaches to the selection of hyper-parameters, and other design decisions in the experimental setup. Next, the experimental section discusses the experimental results and lessons we have learned. We summarize the paper with a recommendation to practitioners and our suggestions for future work.

## II. ANOMALY DETECTION CONTEXTS

While many practitioners are eager to see which method is best for their application, the specifics of the application

<sup>1</sup>The *shallow* techniques corresponds to those we call *classical*. We prefer the later terminology, as models based on random forests are in their essence deep, although they cannot capture semantic structure — a touted feature of deep models

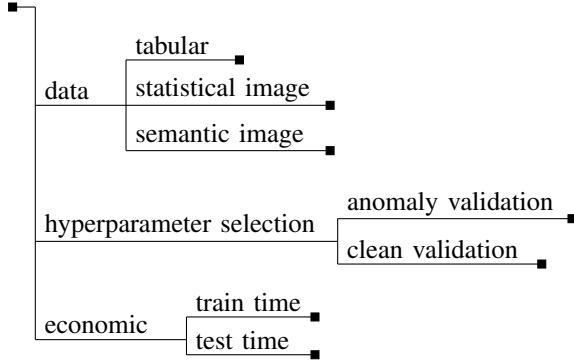


Fig. 1: Various aspects of anomaly detection comparison forming the *context* of the experiment.

may differ. We have conducted a large number of experiments to identify the main sources of variability influencing the performance of anomaly detection methods. The number of combinations of these aspects is huge, therefore we identified the key axes of variability: datasets, hyperparameter selection strategy, and economic point of view. From these axes, we select a few discrete points, on which we will provide a comparison. The particular combination of the selected aspect will be called *context*, see Fig. 1 for illustration.

The first axis is the target data domain. Our experiments used two types of datasets: *tabular* and *image*. This is the most obvious split and indeed most authors of prior art test their methods on either choice of data. Another possible way to look at data is whether they contain *statistical* or *semantic* anomalies. Statistical anomalies should be located in areas of low likelihood of the normal class, while semantic [1] anomalies cannot be differentiated from normal data in a statistical way. This is because they appear in datasets with multiple sources of variations, where only some of them are considered anomalous. Such types of anomalies are most common in image datasets. Imagine a detector that aims to learn a representation of birds from images without preprocessing. Most of the bird pictures are going to have sky in the background. Since the background occupies most of a picture and therefore has a strong signal, a bird on grass is going to be a statistical anomaly, while an airplane with sky in the background is an example of a semantic anomaly in case the original goal was to identify pictures that do not contain a bird. The suitability of the tested methods for the dataset context axis is studied in Sec. V-A.

The second axis of variability is the hyperparameter selection strategy. It should be a gold standard that the experiments are repeated on different splits of data to training, validation, and testing subsets, especially if the datasets are small. However, in most of the reviewed recent papers [55, 89, 79, 3, 67], this procedure was not mentioned with the exception of [76]. Therefore, our comparison fills this gap. Also, it is important to define the nature of information available for selection of the hyperparameters: it is indeed a very different task if there is some (often small) number of known anomalies in the validation dataset that can be used to choose hyperparameters by cross-validation or if the validation dataset is clean. In our

experience, the former case is more common. Our observations are summarised in Sec. V-B.

The third axis is the economic aspect of a problem. There might be serious computational restrictions present in solving real-life problems. One might then not opt for a method that promises state-of-the-art performance, but for another that reaches slightly worse performance but can be trained economically and its performance is robust with regards to hyperparameter optimization. More details on this can be found in Sec. V-C.

Finally, Sec. V-D contains other influences that we have originally considered to be important, but in the end did not prove to make a significant difference in comparison of multiple methods. These include the use of performance measures other than traditional AUC, the use of Bayesian optimization and others.

### III. COMPARED METHODS

This section briefly reviews deep generative models in the order of exactness of calculation of likelihood. Therefore, it starts with flow models, continues with probabilistic (variational) autoencoders, where the prior art on the application in anomaly detection is rich, and finishes with generative adversarial networks where the calculation of any score related to likelihood is dubious at best. We specifically focus on issues that affect the performance of the method for anomaly detection, most often the anomaly score, if it is not rigorously defined. We also briefly review other examples of deep methods that are relevant for comparison such as two-stage models and distance-based models that are not generative but can be used in anomaly detection. We do not review classical methods here, as this has been done many times elsewhere, but we list them in the relevant experimental section.

Before the description, we introduce a notation. Training samples,  $\mathbf{x}$ , are assumed to be i.i.d from the underlying probability distribution  $p(\mathbf{x})$  defined on the input space  $\mathcal{X}$ . Following the conventional definition of an anomaly [7], each anomaly detection method is expected to provide a quantity (called score and denoted  $s(\mathbf{x}')$ ) related to the probability of a sample  $\mathbf{x}'$  being generated from  $p(\mathbf{x})$ . The score does not need to be a normalized distribution, as the threshold is typically determined as an empirical estimate of the quantile. Most functions in this section are assumed to have parameters optimized during training.

#### A. Normalizing flows

The name normalizing flows refers to methods relying on the change of variables formula

$$p(\mathbf{x}) = p(\mathbf{z})|\det J_f(\mathbf{z})|^{-1}, \quad \mathbf{z} = f^{-1}(\mathbf{x}), \quad (1)$$

where  $J_f(\mathbf{z})$  is Jacobi matrix of function  $f$  evaluated at  $\mathbf{z}$ .  $p(\mathbf{z})$  is a known distribution of the latent variable  $\mathbf{z}$  from space  $\mathcal{Z}$  of the same dimension as  $\mathcal{X}$ .

Theoretical reviews [63, 49] require  $f$  to be invertible and both  $f$  and  $f^{-1}$  to be differentiable. Flow models therefore primarily differ in how they define the class of functions  $f$ , which ranges from simple affine transformations to solutions

of ordinary differential equations. The expressive power comes from their composition as is usual in neural networks. In the comparison, we consider flows on tabular data only, for which we have implemented well-known RealNVP [22] and MAF [64] flows alongside with a promising class of Sum Product Transform networks — SPTN [69] combining normalizing flows with a graphical model. The likelihood is used as a natural anomaly score.

Flow models have not yet enjoyed a lot of popularity in anomaly detection [93, 80, 21, 69] in comparison to autoencoders reviewed below. To us, this is surprising, since these methods can exactly calculate likelihood functions, which under a good fit are the ideal anomaly score. Meanwhile, the focus of the surrounding community is on the topic of *out of distribution detection* (OOD)<sup>2</sup> [59], which is very related to anomaly detection if not being equal. Ref. [17] suggests to use ensembles, while [74] recommends to convert the single-class problem to classification problems in the spirit of [85]. A deep investigation of OOD in [47], shows that with low-level features such as pixel intensities, flows tend to learn local models, i.e. according to taxonomy in [75] they fail to detect semantic anomalies.

### B. Autoencoder-based models

Autoencoder-based models differ from flows by relaxing the exact mapping between  $\mathbf{x} = f(\mathbf{z})$  (1) into a probability distribution  $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{z}), \text{diag}(\boldsymbol{\sigma}_{\theta}(\mathbf{z})))$ ,<sup>3</sup> called *decoder*. The symbol  $\theta$  denotes the trainable parameters of the decoder, e.g. weights of a neural network. The marginal likelihood is computed as

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}, \quad (2)$$

where  $p(\mathbf{z})$  is a chosen prior probability distribution of the latent variable. This relaxation allows for more flexible models, e.g. using different dimension of  $\mathbf{x}$  and  $\mathbf{z}$ . However, training and evaluation of the model is more demanding since the marginal likelihood (2) is not available in closed form. Therefore [45] introduces *encoder* distribution  $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}(\mathbf{x})))$  with parameters  $\phi$  allowing approximation of Equation (2) as described below in Equation 3.

Various modifications of the original formulation have been proposed, giving rise to many specialized methods. Below we describe extensions in three blocks according to i) approximation of the likelihood (2) used for training, ii) prior model, iii) approximations used for evaluating the anomaly score, and iv) various modifications of the original concept.

a) *Training loss*: The original Variational Autoencoder [45] (VAE) proposes to replace (2) by the evidence lower bound (ELBO)

$$\mathcal{L}_{\text{VAE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = -\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})), \quad (3)$$

<sup>2</sup>Out of distribution detection means identifying samples coming from a different dataset. For example, a model trained on MNIST / CIFAR10 should assign a low likelihood to samples from Fashion MNIST / SVHN respectively.

<sup>3</sup>Other forms of the distribution are possible, e.g. Bernoulli for scaled pixel intensities.

which combines reconstruction error with regularization term in form of the Kullback-Leibler divergence (KLD) between the encoder distribution and the prior. Models based on (3) will be referred to as the VAE family.

Asymmetry of the KL divergence motivated search for a more accurate metric. Ref. [86] proposes to replace KL by a Wasserstein divergence, yielding training loss function in the form:

$$\mathcal{L}_{\text{WAE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = -\mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + \lambda D(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})), \quad (4)$$

where  $\lambda > 0$  is a scalar hyperparameter, an  $D$  is an arbitrary divergence. The most commonly used divergence is the kernelized maximum-mean-discrepancy (MMD) with kernel  $k$  which was reported to performs well in matching high dimensional distributions [96]. Models based on (4) will be referred to as the WAE family.

An alternative choice of the divergence  $D$  in (4) proposed in [86] is the adversarial loss, which in combination with the Gaussian decoder coincides with the adversarial autoencoder [56]. This divergence introduces a third network  $d_{\psi}(\mathbf{z}) : \mathcal{Z} \rightarrow [0, 1]$ , called discriminator, trained to distinguish between samples from the prior  $p(\mathbf{z})$  and samples  $\mathbf{x}$  projected by the encoder  $q(\mathbf{z}|\mathbf{x})$ . Every step of optimization separately updates the autoencoder and discriminator parts to minimize the loss functions

$$\begin{aligned} \mathcal{L}_{\text{AE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) &= -\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \lambda \log d_{\psi}(\mathbf{z}^q), \quad (5) \\ \mathcal{L}_{\text{D}}(\boldsymbol{\psi}) &= \log d_{\psi}(\mathbf{z}^p) + \log(1 - d_{\psi}(\mathbf{z}^q)), \quad (6) \end{aligned}$$

respectively, where  $\mathbf{z}^p \sim p(\mathbf{z})$ ,  $\mathbf{z}^q \sim q_{\phi}(\mathbf{z}|\mathbf{x})$ . Models trained with loss function (5) will be denoted as AAE.

b) *Prior model*: A common criticism of the VAE model is its use of the standard Gaussian prior  $p(\mathbf{z})$ , which stimulates the distribution  $q(\mathbf{z}|\mathbf{x})p(\mathbf{x})$  to have a single mode, and therefore it is hard to fit data with a multi-modal latent distribution. Ref. [87] proposes a learnable multimodal *Vamp* prior realized as a mixture of  $K$  independent Gaussian components. Vamp prior is compatible with AAE and WAE models since it does not have an analytical expression of KLD in (3). The mean values of components of the mixture are learned together with the parameters of the autoencoder. In the model selection below, the Vamp prior is considered as a binary hyperparameter with an additional parameter,  $K$ , specifying the number of components.

c) *Anomaly Score*: The likelihood function (2) also constitutes the ideal anomaly score. Some training losses such as ELBO (3) were designed as approximations of the likelihood and can thus be used as anomaly scores. However, this interpretation is not so clear for other training losses, i.e. (4), (5), hence their authors propose anomaly scores as part of the method. Nevertheless, many scores are interchangeable, giving rise to another degree of freedom (hyperparameter) for the use of autoencoders in anomaly detection. A common score is based on the first term in the loss i.e. a Monte Carlo estimate of the expectation of conditional log-likelihood over

the encoder, yielding

$$s_{\text{rs}}(\mathbf{x}) = -\frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x} | \mathbf{z}_l), \mathbf{z}_l \sim q_{\phi}(\mathbf{z} | \mathbf{x}). \quad (7)$$

This score, called sampled reconstruction error (abbreviated as rs), was shown in [92] to be more accurate than evaluating (2) by sampling  $z$  from the prior  $p(z)$ . Further simplification is based on replacing samples from the encoder by its mean, yielding the common reconstruction error score (abbreviated as rm)

$$s_{\text{rm}}(\mathbf{x}) = -\log p_{\theta}(\mathbf{x} | \boldsymbol{\mu}_{\phi}(\mathbf{x})) \quad (8)$$

The usage of (8) is justified by the assumption that taking the mean at the encoder should approximate (7) while having lower computational demands. For adversarial autoencoders, these simplifications can be combined with the discriminator score [79, 95],

$$s_{\text{a}}(\mathbf{x}) = \alpha s_{\text{rm}}(\mathbf{x}) + (1 - \alpha) d_{\psi}(\boldsymbol{\mu}_{\phi}(\mathbf{x})), \alpha \in [0, 1]. \quad (9)$$

The reconstruction-error based anomaly scores were criticized in [70] for not capturing the true data density  $p(\mathbf{x})$ . The proposed replacement is based on the orthogonal decomposition of the data into  $x = x^{\perp} + x^{\parallel}$  where the  $x^{\parallel}$  lies in the tangent space of to the manifold defined by the decoder. This allows to decompose the marginal likelihood into a product of two orthogonal parts

$$p(\mathbf{x}) \approx p(\mathbf{x}^{\parallel}) p(\mathbf{x}^{\perp}), \quad (10)$$

where  $p(\mathbf{x}^{\perp})$  is the reconstruction error, and  $p(\mathbf{x}^{\parallel})$  is obtained by transformation of variables (1). This score is abbreviated as  $jc$  in the following text. The calculation of (10) is expensive, as it needs to compute the singular value decomposition of the Jacobian. For implementation details, see [70] or [84].

d) *Other models and techniques:* A plethora of models based on probabilistic autoencoders and specialized for anomaly detection was introduced in recent years, such as [98, 66, 92, 72, 16, 15]. Below, we list models included in the comparison and not described above.

The self-adversarial Variational Autoencoder (adVAE) [89] was included because it claims superiority over the state-of-the-art methods, such as VAE, DAGMM [98], WGAN-GP [36] or MO-GAAL [55] on tabular datasets. It augments the usual encoder-decoder pair with a transformer, whose goal is to simulate anomalies during training. The seeming flaw of the model is that it is trained only on normal data and there is no link between real and simulated anomalies. The sampled reconstruction is used as an anomaly score.

Despite its name, GANomaly [3, 2] is more related to adversarial autoencoders than to GANs. It consists of encoder-decoder-encoder with a discriminator, similar to an AAE. The anomaly score is the difference between latent representations of a sample after the first and second encoding. An upgrade to this model, skip-GANomaly [4], uses skip connections in a U-Net type architecture. Here, the anomaly score is a combination of reconstruction error and feature-matching loss (see the next section on fmGAN). Although originally proposed only for use in images, we have implemented a variant for tabular data as well.

### C. Generative adversarial networks

Generative adversarial networks [33] construct and train two networks: a generator  $g_{\phi}(z) : \mathcal{Z} \rightarrow \mathcal{X}$  and a discriminator  $d_{\psi}(x) : \mathcal{X} \rightarrow [0, 1]$  approximating the probability of  $x$  being a sample from the data distribution rather than the generator. The generator aims to transform samples from  $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  to  $\mathcal{X}$  such that they are indistinguishable from the real data. Formally, the optimization objectives are

$$\mathcal{L}_d(\psi) = \log d_{\psi}(x) + \log(1 - d_{\psi}(g_{\phi}(z))), \quad (11)$$

$$\mathcal{L}_g(\phi) = -\log d_{\psi}(g_{\phi}(z)), \quad (12)$$

where  $\mathcal{L}_d$  is maximized while  $\mathcal{L}_g$  is minimized,  $z$  is sampled from the prior and  $x$  from the training set. The optimization searches the saddle point of the two losses, which is difficult and notoriously unstable. Therefore a long series of work, e.g. [40], proposes improvements over the basic approach [33]. One of the approaches is based on the introduction of feature-matching loss [77]. We will denote the model trained with this loss as feature-matching GAN (fmGAN). In fmGAN training, the cost function of the generator is augmented with output of some intermediate layer of the discriminator. Specifically, the generator is optimized as

$$\mathcal{L}_{\text{fm}}(\phi) = \alpha \mathcal{L}_g(\phi) + \|h_{\psi}(x) - h_{\psi}(g_{\phi}(z))\|^2, \quad (13)$$

where  $h_{\psi}$  is the output of the intermediate layer of the discriminator and  $z$  is a sample from  $p(z)$ . This feature-matching loss is used in AnoGAN [79] for detection of anomalous objects in images, with hyperparameter  $\alpha$ , which was zero in the original publication [77].

GANs are frequently augmented with a third model  $q(z|x)$  [24] which makes the distinction between GANs and VAEs blurred, as demonstrated by using min-max (GAN-like) approximation of Wasserstein divergence in Wasserstein autoencoders, Sec. III-B. This makes it sometimes hard to assign a model to some class (for example GANomaly belongs according to us to probabilistic autoencoders).

Recall the role of the discriminator is to discriminate *generated* samples from *real* ones. Since the generator is trained to generate samples with a high discriminator score, it seems logical to use the discriminator to score anomalies

$$s_{\text{GAN}}(\mathbf{x}) = 1 - d_{\psi}(\mathbf{x}), \quad (14)$$

which is used e.g. in [55]. The common critique is that the discriminator was not trained to recognize an arbitrary distribution of the anomalies, but only that of the latent transformed by the generator. Thus it may fail to recognize anomalous samples of interest. AnoGAN recognizes this flaw and proposes to augment the discriminator loss (13) for training and an iterative procedure that searches for the latent code  $z$  most likely to generate the tested sample to identify anomalous images. However, this procedure is computationally expensive. Its sequel, f(ast)AnoGAN [78], uses Wasserstein GAN with gradient penalization [36] to improve stability and adds  $q(z|x)$  to find  $z$  closest to given  $x$  in  $q(x|z)$  faster. The anomaly score of fAnoGAN is a combination of discriminator score and feature-matching loss.

Multiple-Objective Generative Adversarial Active Learning (MOGAAL) [55], train  $k$  generators against a single discriminator on input data divided into  $k$  subsets. The usual discriminator score in Equation (14) is used to test new samples. Other anomaly detection models derived from GAN include [95, 48, 67], however, it seems that autoencoder-based models are more popular for anomaly detection and our selected candidates are representative.

#### D. Two-stage models

A recurring idea [27, 94, 76, 83] is to combine autoencoders with a secondary model acting on the latent space defined by the encoder. The motivation behind it is that the encoder should preserve the semantic information of the sample and remove noise (e.g. background in images). Additionally, reducing the size mitigates the curse of dimensionality, as high dimensions can be problematic for some models.

We are not aware of a general term for this approach. We use the term *two-stage models*, following [18], although [13] uses the term *deep hybrid models*. In [76], the model optimizes the projection of data (by virtue of NNs) to a new space, where they can be easily enclosed in a sphere of minimum radius. The approach presented in [83, 94] explicitly splits the creation of the detector into two parts. It first trains a VAE (and its variants) and then it fixes the encoder. The anomaly score is calculated by a kNN [83] or by OC-SVM [94] detectors in the latent space, obtained by projecting the sample by the fixed encoder. The two-stage models can be also viewed as a kNN with a trained metric or OC-SVM with a trained kernel. The embedding can be optimized differently, for example by enforcing the margin between anomaly candidates and normal data as done in the REPEN [61] method, which uses an ensemble of 1NN detectors as the second stage.

## IV. EXPERIMENTAL SETUP

### A. Datasets

The choice of datasets (mainly the tabular ones), was guided by two criteria: first, they ought to be publicly available and second, they should appear in surveys or in articles presenting new methods. In total, we have collected 40 tabular datasets, the majority of which came from the UCI repository [25]. With the exception of ANNThyroid, Arrhythmia, HAR, HTRU2, KDD Cup 99 (small), Spambase, Mammography, and Seismic, where the anomaly class has a clear meaning (security incident or disease), we have followed the technique of [26] for creating artificial datasets for anomaly detection tasks from classification datasets. More precisely, we have used only "easy" and "medium" anomalies, as "hard" and "very hard" are not truly anomalous in the sense of being clearly statistically distinct from the normal class. Prior to model training, features on tabular datasets were normalized to have zero mean and unit variance. Further details of the datasets are provided in Tab. I.

The number of image datasets used for the evaluation of deep models is limited, as there are very few datasets designed purely for anomaly detection - MNIST-C [58] and MVTEC-AD [8]. Therefore, we have decided to extend these with artificially created anomaly datasets based on common image

datasets MNIST [53], FashionMNIST [91], CIFAR10 [51], and SVHN2 [60]. These are also used for anomaly detection tasks in the prior art [67, 70, 76]. Most often a single class (of digits/objects) is considered normal and the rest anomalous, which is the primary setting of our experiments. Only three subsets *wood*, *grid*, *transistor* of the MVTEC-AD dataset were used due to our computational constraints. These were chosen since they represent problems of various degrees of difficulty. Downsampling to  $128 \times 128$  was required to fit the computational envelope for most methods on the real-world high-resolution images in MVTEC-AD. We have linearly extrapolated the image data when training GANomaly so that the input dimensions were a multiple of 16 (this is a result of the model's fixed architecture). No other preprocessing has been applied prior to training since the source data already have all channels scaled to  $[0,1]$ . Basic statistics on image datasets are shown in table Tab. II.

We have performed a visual inspection of the nature of anomalies in the image datasets, see Supplementary C. Since most of the datasets that were manually processed (MNIST, FashionMNIST, MVTEC-AD and MNISTC) have a rich and consistent number of samples in the normal class and clear anomalies, we consider them to be statistical anomalies. On the other hand, images in the majority of classes in CIFAR10 and SVHN2 have strong background and are thus considered to contain semantic anomalies. This prior division is also supported by a different behavior of different methods as reported in the Supplementary, Figure D.1.

### B. Data splits and experiment repetitions

The experiments on tabular data and MNIST-C/MVTEC-AD image datasets were repeated five times with different random cross-validation splits. Specifically, in each repetition (five in total) of an experiment with the same model hyperparameters, the *normal data* in each dataset was randomly split in 60%/20%/20% ratios to train/validation/test subsets, respectively. *Anomalous data* were split such that 50% were in the validation part and 50% in the testing part, which means the training subset has not contained anomalous samples.<sup>4</sup> The proportion of anomalies that were used in the validation phase varied from zero to the selected 50%.

On the rest of the image datasets, due to the already substantial computational requirements, we have not trained models with the same hyperparameters on repeated random cross-validation splits. In our experience, the results on different splits of these datasets are almost the same, since the number of samples is large and our trial experiments (see Supplementary Tab. A.1) have not exhibited a large variation between the different random cross-validation experiment repetitions.

### C. Notes on implementation of models

As mentioned in the introduction, we have compared various types of *deep* methods to *classical* ones serving as an etalon. Classical methods included ABOD [50], HBOS [31],

<sup>4</sup>A training set without any anomalies is in practice very optimistic, but this decision removes another degree of freedom from the evaluation for the sake of clarity of results.

dataset	alias	dim	anom	normal
ANNthyroid	ann	21	534	6665
Arrhythmia	arr	275	206	245
HAR	har	561	1944	8355
HTRU2	htr	8	1638	16257
KDD99 (10%)	kdd	118	396742	97276
Mammography	mam	6	260	10921
Seismic	sei	24	170	2412
Spambase	spm	57	1812	2786
Abalone	aba	10	50	2151
Blood Transfusion	blt	4	16	382
Breast Cancer Wisconsin	bcw	30	206	356
Breast Tissue	bts	9	22	65
Cardiotocography	crd	27	228	1830
Ecoli	eco	7	108	205
Glass	gls	10	94	112
Haberman	hab	3	14	225
Ionosphere	ion	33	122	225
Iris	irs	4	46	100
Isolet	iso	617	3300	4496
Letter Recognition	ltr	617	3600	4196
Libras	lbr	90	142	215
Magic Telescope	mgc	10	3882	12331
Miniboone	mnb	50	23922	93565
Multiple Features	mlt	649	800	1200
PageBlocks	pgb	10	384	4911
Parkinsons	prk	22	44	146
Pendigits	pen	16	5384	5537
Pima Indians	pim	8	176	500
Sonar	snr	60	96	110
Spect Heart	sph	44	52	211
Statlog Satimage	sat	36	2630	3592
Statlog Segment	seg	18	938	1320
Statlog Shuttle	sht	8	28	57767
Statlog Vehicle	vhc	18	132	627
Synthetic Control Chart	scs	60	200	400
Wall Following Robot	wrb	24	2220	2921
Waveform-1	wf1	21	1482	3302
Waveform-2	wf2	21	1472	3302
Wine	wne	13	70	106
Yeast	yst	8	390	751

TABLE I: Basic statistics of the tabular dataset designed for anomaly detection (above split) and multi-class datasets (below split).

dataset	alias	dim	anom	normal
MNIST-C	mnisc	28x28x1	70000	70000
MVTec-AD - wood	wood	1024x1024x3	60	266
MVTec-AD - grid	grid	1024x1024x3	57	285
MVTec-AD - transistor	transistor	1024x1024x3	40	273
CIFAR10	cifar10	32x32x3	54000	6000
FashionMNIST	fmnist	28x28x1	63000	7000
MNIST	mnist	28x28x1	63686	6312
SVHN2	svhn2	32x32x3	80327	18960

TABLE II: Basic statistics of image datasets designed for anomaly detection (above split) and multi-class datasets (below split).

class	model	acronym	class	model	acronym
flows	MAF	maf	two-stage	DAGMM	dgmm
	RealNVP	rnvp		DeepSVDD	dsvd
	SPTN	sptn		REPEN	rpn
		VAE-kNN		vae	
autoencoders	AAE	aae	VAE-OC-SVM	vao	
	adVAE	avae			
	GANomaly	gano	ABOD	abod	
	skipGANomaly	skip	HBOS	hbos	
GANs	VAE	vae	classical	IsolationForest	if
	WAE	wae		kNN	knn
	fAnoGAN	fano		LODA	loda
	fmGAN	fmgan		LOF	lof
	GAN	gan		OC-SVM	osvm
	MOGAAL	mgal	PidForest	pidf	

TABLE III: Overview of the main classes of compared methods and the acronyms used in the text.

LODA [68], LOF [11], IsolationForest [54], OC-SVM [81], PIDForest [34], and kNN [73]. For ABOD, HBOS, and LODA, we have used pyOD library [97] implementation, for LOF, IsolationForest, and OC-SVM we used scikit-learn [65] implementation, and last but not least we have used our own implementation of kNN. The acronyms used in the result section are summarized in Tab. III together with the classification of the deep methods as described in Sec. III.

Since image datasets are typically much larger than tabular, the OC-SVM was implemented as an ensemble of 10 OC-SVM models trained disjoint subsets of data, see Sec. F, because it has at best  $O(n^2)$  scaling in the number of samples.

To ensure consistency among deep models, we implemented all methods except the MOGAAL<sup>5</sup> ourselves using the Flux [41] framework in Julia [10] language. Apart from the models mentioned in Sec. III, we have also implemented DeepSVDD [76], DAGMM [98] and REPEN [61], which have been included in multiple comparisons [75, 89, 14].

We emphasize that while implementing models, we have carefully compared our implementations to the reference where possible and (or) verified that our experimental results are similar to those provided in the corresponding publication.

Neural networks were trained with the ADAM [43] optimizer with early stopping measuring the continued decrease of loss on the validation dataset. All deep models trained on all image datasets used convolutional layers. More implementation details are in the Supplementary materials Sec. F. The implementation code in the form of a Julia package can be found at <https://github.com/aicenter/GenerativeAD.jl>.

#### D. Hyperparameters and their optimization

Properly exploring the space of hyperparameters of all models is paramount to achieving fair and comparable experimental comparison, yet this is often superficially treated. Researchers often use *default* or *recommended* values ignoring that they are sub-optimal on datasets they use in their comparison. A nice demonstration of this are conflicting results of the MOGAAL method in the original publication [55] and

<sup>5</sup>The pyOD implementation was used.

in [89]. A prototypical example in the classical methods is OC-SVM, which is typically used with Gaussian kernel and with  $\nu$  set to some default value, for example 0.05 [68], but can achieve better results with different kernels. The choice of hyperparameters in anomaly detection is everything but easy. But this means that the experimental settings should be set up such that all methods have been optimized equally. We conjecture that recommended and default values of hyperparameters are strongly correlated with the choice of evaluation datasets in the publications that recommend them.

*Random grid search:* In order to explore the hyperparameter space of each method properly, we have employed a random search over a predefined grid for each method. This allows the construction of sections through the space for sensitivity studies. Moreover, it is frequently more efficient than grid search [9] and more flexible. For each model, dataset, and repetition, we have sampled 100 configurations from corresponding sets (see Tab. A.3–A.5) and trained the models with them. In order to keep the neural-network-based models fixed across the repetitions on a single dataset, for each hyperparameter configuration we have also sampled a random seed that was used to initialize the network weights. To prevent running the training of expensive methods forever, there was a hard deadline of 24 hours in which the training of a single configuration for a given number of repetitions/classes should be finished. This automatically penalizes complicated models.

Encoders for the two-stage models were selected from models performing best in terms of validation AUC or reconstruction error on the validation set. The second-stage models (kNN and OC-SVM) used hyperparameters sampled from Tab. A.3.

*Bayesian optimization:* To overcome the limitation of random sampling in larger configuration spaces we have also trained models whose hyperparameter choice was guided by Bayesian optimization of an evaluation metric on validation data. We followed 50/50 strategy, where the underlying Gaussian process has been fitted with 50 randomly sampled configurations and another 50 have been sampled based on the acquisition function. We have relied on the off-the-shelf implementation from the scikit-optimize framework [39] with default parameters.

*Number of anomalies in the validation set:* A thorough exploration of the hyperparameter selection context also requires changing the criteria of model selection. When anomalies are available for validation, we select hyperparameters maximizing the AUC on the validation set. For experiments with no available anomalies, we have decided on the following hyperparameter selection mechanism. For *classical* methods, we have used default hyperparameter values from literature - either they were recommended by authors of the method, were used in a survey, or are default in a given implementation. Their overview is in Tab. A.6. For *deep* methods, this is unfortunately impossible, since their hyperparameter space is much larger and the values are usually tuned to a specific dataset. Therefore, in order to have a universal solution, we have selected the already trained and evaluated models based on the lowest anomaly score on validation normal data. This approach is theoretically justified for models with proper

likelihood.

*Ensembles:* Some results on ensembles of anomaly detectors were already reported in [17, 59]. Since the other experiments required training a large number of models, it was decided to test whether even a naive approach to this problem brings some improvements. In order to mitigate some uncertainty given by different hyperparameter values, ensembles of detectors were constructed by averaging scores of some number of best-performing detectors. We have experimented with different fixed sizes of ensembles — either top 10 or top 5. In such a case, the anomaly score is the average of the anomaly scores of ensemble members.

*Performance criteria:* While the area under the ROC curve (AUC) is the most common criteria in anomaly detection, some authors use different metrics, such as partial AUC [23] or the true positive rate at a chosen false negative rate (TPR@). We have re-evaluated all results obtained for the random grid search on the TPR@5%.

We have kept track of the time spent on training of individual models and also on the time needed to evaluate them on validation and test sets. In total, we have trained 1,364,989 model instances in 9619 CPU days, evaluated 4,256,470 different scoring functions in 2704 CPU days, and created 10.2TB of experimental data.

## V. EXPERIMENTAL RESULTS

Before starting with the description of the experimental results, here we summarize conventions that are used unless said otherwise. The performance results are estimates of the AUC on the testing set and are averaged over five random cross-validation repetitions in the case of tabular and MvTec-AD/MNISTC image datasets. When ranks are reported, they are calculated by ordering methods on each dataset and calculating the average across them (as recommended in [20]). Hyperparameters are selected using the best average performance over 5 seeds, or individually over 10 anomaly classes because the individual class splits constitute different anomaly detection problems.

### A. Dataset context

The results of experimental comparison on all dataset types are presented in the form of critical difference diagrams (CDD) as recommended by Demšar [20], are in Fig. 2. Diagrams show the average rank of detectors across the datasets together with a confidence band that indicates that a statistical test cannot reject the hypothesis that two detectors perform the same. The underlying AUC values on the testing set for all individual datasets are given in Tab. B.1, B.2, B.3, and B.4. We now comment on the influence of the datatype with respect to two types of hyper-parameter selection strategies differing in the number of anomalies in the validation set as defined in Section II: i) anomaly validation context, and ii) clean validation context.

*Tabular data:* OC-SVM works the best and it is *statistically better* than almost all detectors except autoencoder-based generative models and VAE combined with OC-SVM in the case of anomaly validation context. The first 11 places (roughly one

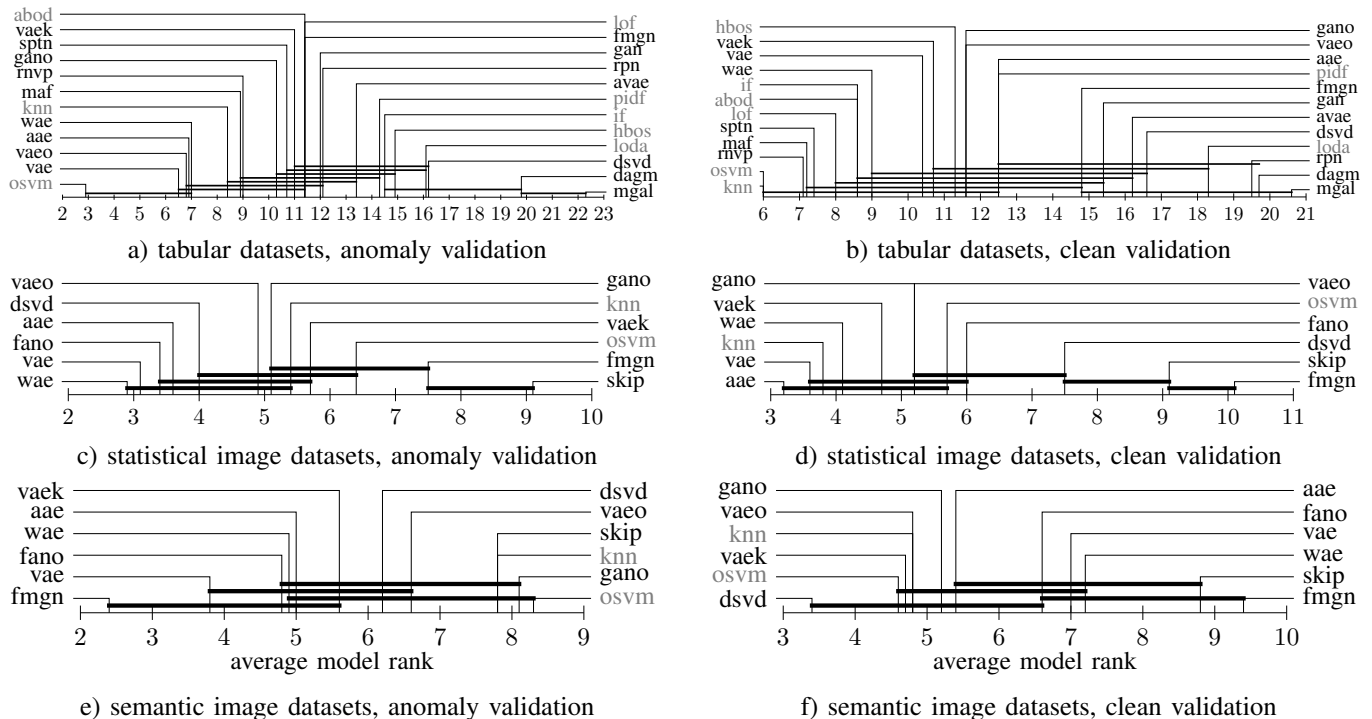


Fig. 2: Critical difference diagram of models ranked via the test AUC. Models whose performance is statistically indistinguishable have difference of ranks under the critical value of the Nemenyi test  $CD_{0.1}$  and are joined by a horizontal band. Results are presented for different types of datasets: tabular (Top row), image datasets with statistical anomalies (Middle row), and image datasets with semantic anomalies (Bottom row); and two different hyperparameter selection cases: using anomalies in validation (left) and using clean validation (right).

half) belong to models that can be divided into three groups: (i) OC-SVM and its variants, which estimate a density level of a distribution; (ii) flow models and kNN which estimate the pdf (un-normalized in case of kNN); (iii) and variants of auto-encoders, where reconstruction error is related to pdf as explained in [84]. The same types of methods occupy the top positions in the clean validation context, Fig. 2b), however in a different order. The best is the kNN and all other pdf-modelling methods (flows) have improved relative to the anomaly validation context. The autoencoder-based methods moved beyond classical methods (LOF, ABOD, IF). We believe that models that are in the lower half of the scale in both validation contexts are not suitable for detection of statistical anomalies. We cannot explain the poor performance of MOGAAL, DAGMM and adVAE<sup>6</sup> and we attribute it to different experimental environment. DeepSVDD was primarily implemented for image problems, where it performs relatively.

Moreover, differences in mean ranks of many models in Fig. 2 are statistically insignificant at level  $p = 0.1$ , which is disappointing. Assuming the ranks remain the same, another 51 datasets would be needed to make the difference between OC-SVM and VAE statistically significant on tabular data with 50% anomalies. This indicates that the results are still noisy and can be easily changed for a different choice of datasets.

<sup>6</sup>We have contacted the author of pyOD from wherein we took the implementation of MOGAAL and we were assured that his implementation is a copy of that provided by the authors. Therefore, we consider the implementation to be correct.

*Statistical image data:* WAE and VAE models have the best average rank when evaluated on statistical image data, although their lead is not statistically significant over most of the other models as is evident from Fig. 2c). The autoencoder-based methods (AAE,VAE,WAE) perform well also in the clean validation context, complemented by kNN, Fig. 2d).

*Semantic image data* A different story is being told by Fig. 2e) where the ranking of methods on image datasets with semantic anomalies is dominated by fmGAN by a large margin in the anomaly validation context. However, it is also the worst method in the clean validation context. In an opposite manner, OC-SVM and kNN perform very poorly in the anomaly validation context, but they are among the best in the clean validation context. The best performing method in the clean validation context is DeepSVDD [76]. We conjecture that with increasing number of anomalies in the validation set, the problem is approaching that of supervised classification with hyper-parameters playing the role of parameters.

The typical anomalies detected by various methods on image datasets are provided in the Supplementary C.

### B. Hyperparameter selection context

The influence of hyperparameter selection procedure on the results in the previous section is now studied in detail for few selected methods. We choose only those that scored among the best in the previous Section. First, we analyze the sensitivity of these methods to the number of anomalies in the validation set.



Second, we study hyperparameter selection for two individual methods, variational autoencoder family and OC-SVM.

1) *Impact of the number of anomalies in the validation set:* The process of hyperparameter selection described in Sec. IV-D depends on the availability of examples of anomalies in the validation set (recall that it is assumed that the validation dataset does not contain unknown anomalous samples, i.e. is not contaminated). Fig. 3 displays the influence of the number of anomalous samples in the validation set on a finer grid between the two contexts reported before. Note that for the first point on the x-axis, *clean*, the mechanism of model selection was different from the rest of the graph, as noted in Sec. IV-D. This is the reason for the significant difference between the clean context and the remaining points.

First, we observe that the quality of the models selected using anomalies improves with an increasing number of anomalous samples which is expected. However, for low number of anomalies, many methods perform significantly worse than in the case of clean validation context. This behavior is notable across dataset types, especially for OC-SVM, and to some extent VAE. We conjecture that the hyper-parameter selection procedure of those methods has a tendency to overfit and its hyperparameters are not robust. In contrast to this, the performance of kNN, WAE and RealNVP degrades slowly with declining number of anomalies, which suggests that they are quite robust in difficult operating conditions. We attribute it to the fact that these methods are more exact in their estimation of data likelihood than the rest.

Second, we notice that the experimental results on the semantic image datasets are generally poor, as the AUC of the best model (fmGAN) on CIFAR10 is 0.72 and similarly on SVHN2, where the best model achieved 0.74. On the other hand, anomaly detection methods perform well on statistical image datasets. This indicates, contrary to popular belief, that the models fail to learn or identify the important semantic information, or they consider different semantic information anomalous and they should be told *which* semantic aspect of an image should be considered as an anomaly, as for example blurred images might be anomalous as well.

Results of the same study for individual image datasets are presented in the Supplementary, Table D.1. We also provide an illustration of what images were identified as normal and anomalous for the tested methods, Supplementary C.

A practitioner might also desire a method robust with respect to poor choice of hyper-parameters. In general, deep methods in our experiments have demonstrated higher variance, probably due to the large number of hyperparameters and stochasticity involved in their initialization and training via batched gradient optimization. In this respect, GAN-based models seem to be the least robust, which is in line with [19] stating that GANs are not directly optimized for anomaly detection. This hints at the potential cost of hyperparameter optimization — with higher performance variance, one is less likely to train a well-performing model in a given number of attempts.

2) *Sensitivity Analysis of the VAE family:* Autoencoder-based methods form a whole family with multiple sources of variability, as identified in Sec. III-B to be: i) approximation

of the likelihood in training (loss function), ii) the richness of latent prior, and iii) the anomaly score. We will analyze the sensitivity of the results to these choices on tabular data in the anomaly validation context. We focus on this family since most of the novel deep generative models for anomaly detection are based on the autoencoder architecture. Additional degrees of freedom include the parametrization of variance of  $p_{\theta}(x|z)$ , which could be either fixed (called VAE-constant), used in [94, 89, 2], scalar (called VAE-scalar), or full diagonal (called VAE-diagonal), used in [5, 92, 95]. In the experiments, all three variations were tested on tabular data, however on image data, the full diagonal was skipped due to computational constraints (and in line with the prior art, where only fixed variance is used).

The overall comparison in Fig. 2 revealed that WAE and vanilla VAE variants perform best. The other degrees of freedom, namely richness of prior, used anomaly score and parametrization of variance were treated as hyperparameters. Fig. 4 extends the study by showing the distribution of ranks over tabular datasets for different variants of VAE including GANomaly and adVAE.

First, notice that the spread of the method’s ranks over various datasets is significant, as even ranks of the best methods vary from 3 to 15. This means that the conclusions below need to be taken with a grain of salt, as the experimental results are extremely noisy.

The ELBO-based score, -el, together with the orthogonal decomposition of the likelihood [70], -jc, does not perform well. The sampled reconstruction error (an MC estimate of (7)), -rs, almost always performs better than the usual reconstruction error, -rm, calculated according to (8). This demonstrates the common approach of replacing the mean of the decoder with that of the encoder is inferior but computationally cheaper (see Tab. IV with prediction times). The discriminator score (14), -di, of AAE (an autoencoder combined with GAN) seems to be also on par with the MC estimate (7).

From the same figure, we also conclude that the models modelling full diagonal in  $p_{\theta}(x|z)$ , -d-, seem to be better than the scalar, -s-, or constant, -c-, variants. This result is important, as many comparisons in the prior art use the VAE-constant, despite the version with full diagonal being discussed in the original publication [45].

The rich prior distribution on the latent space proposed in [87], VAMP, -v-, does not seem to give an advantage in the anomaly detection except in the AAE. Similarly, recent variants adVAE and GANomaly do not seem to work well on the tabular data, but they were not evaluated on them in the original publications.

3) *Sensitivity study of OC-SVM:* Domination of OC-SVM on tabular data in anomaly validation context contrasts to many prior experimental comparisons [32, 15, 19, 34, 42, 89]. The search for the culprit found it to be the hyperparameter selection. This study has varied the  $\nu$  parameter, kernels, and their parameters, which is much more than the most prior art does, which is fixing the kernel to RBF and tests few values of its width  $\gamma$  and  $\nu$ . Inclusion of other kernels into the search for hyperparameters seems to be the major source of improvement in this case. Replacing the OC-SVM with

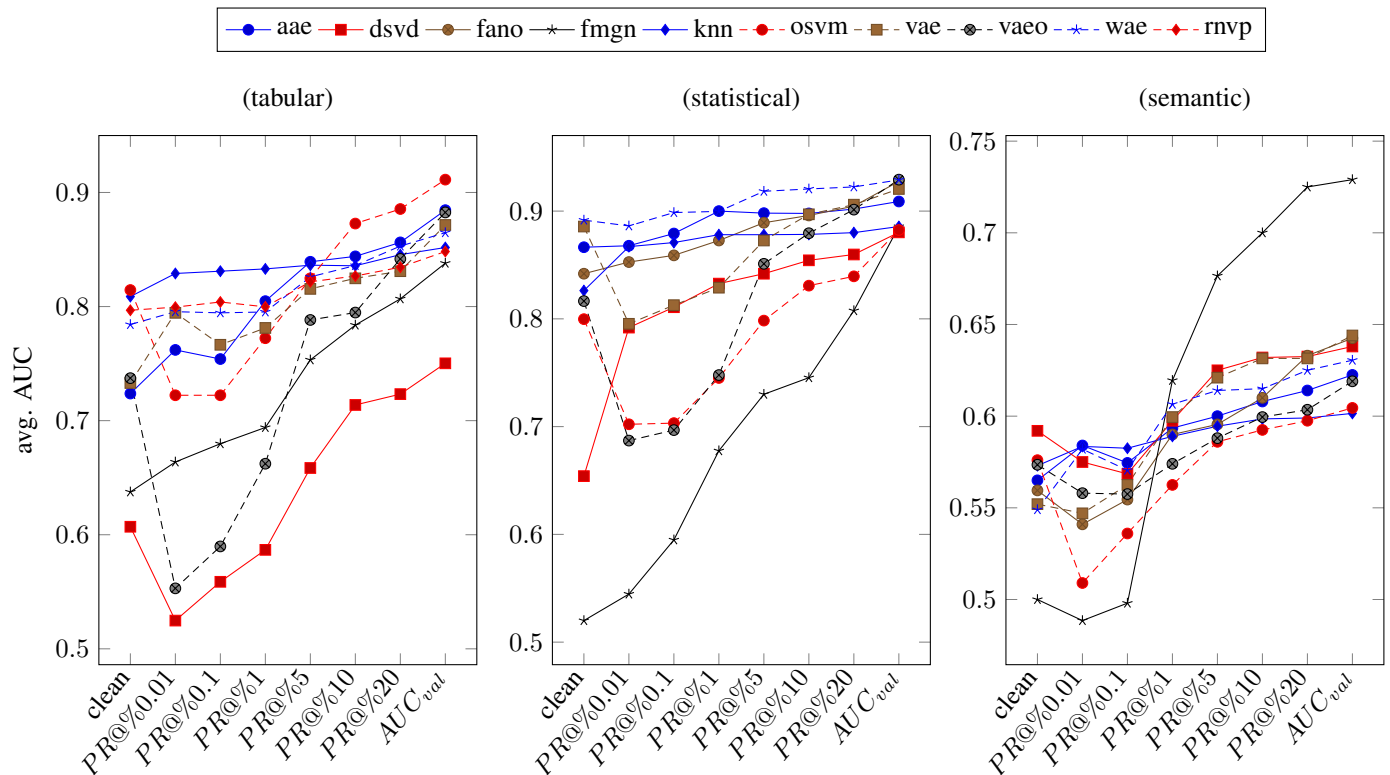


Fig. 3: Sensitivity of methods to the number of anomalies available in the validation set for hyperparameter selection visualized in terms of the achieved AUC aggregated over all datasets in each category (columns). The clean validation context is the left-most point on the x-axis, and the anomaly validation context (50% of available anomalies) is the right-most point. The points in-between were obtained by selecting models with highest precision on the reported portion (e.g. 5%) of validation samples with the highest anomaly scores.

	vae-s-rs	vae-d-rs	vae-s-rm	vae-d-rm	vae-d-jc
$\bar{t}_{pred}$ [s]	12.10	18.51	0.11	0.15	57.31

TABLE IV: Average prediction times on the tabular datasets for different combinations of VAE scores and decoder variance estimations. The -d- part stands for model with an estimate of the full diagonal of the decoder variance, -s- is a scalar estimate. Sampled reconstruction error (7) is denoted as -rs-, -rm is the anomaly score (8) and -jc is (10).

one restricted to use only the RBF kernel and  $\nu = 0.5$  yields an increase in average rank from 2.9 to 8.1 with an average decrease in performance by 0.06, measured in AUC in the anomaly validation context. This version of OC-SVM is then easily surpassed by variational autoencoders and kNN, as demonstrated in Supplementary, Table D.1. The importance of the choice of kernel is furthermore illustrated by the fact that the sigmoid kernel was the optimal choice for 23 datasets, while the RBF kernel only for 13. Ref. [32] mentions that setting  $\nu = 0.5$  provides universally good results, which may be the reason why many authors do not tune it. In theory, it should be set to much lower values ( $\nu = 0.05$ ) corresponding to the presumed low ratio of anomalies in data, but with Bayesian optimization, we found that the best estimate of  $\nu$  was in some cases even higher, such as  $\sim 0.75$  on the statlog-

vehicle dataset.

### C. Economic context

Practitioners ask for fast and accurate algorithms, but these two features rarely go hand in hand, and a decision on a trade-off has to be made. Interesting methods lie on the Pareto frontier, as in absence of external factor, rationally behaving practitioner does not have a motivation to choose a different model.

Fig. 5a-left shows the trade-off between accuracy and training time for tabular data, where the absolute numbers were replaced by average ranks for robustness. The Pareto frontier contains two methods, which are OC-SVM and kNN. The position of OC-SVM is rather surprising, as its training time is known to scale poorly (quadratically) with respect to the number of samples, but it is caused by most of the tabular datasets being small. Different results may arise for a dataset with many data records. Fig. 5a-right shows a similar trade-off between accuracy and testing (inference) time. OC-SVM is still on the Pareto frontier, but it is expensive, as the complexity grows linearly in the number of samples. fmGAN, GAN and DAGMM methods are there as well – these methods have fast inference but lower accuracy.

We provide results averaged over the studied contexts on image data. Due to variability of the results in each context,

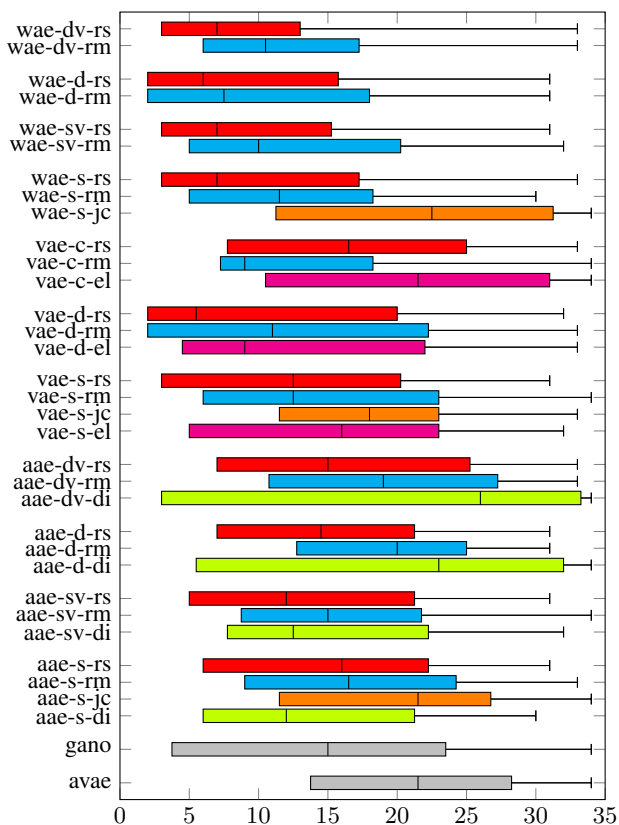


Fig. 4: Sensitivity study of various variants of autoencoder-based methods displayed in the form of boxplots of their ranks in the AUC metric achieved on the tabular datasets. The first three letters of the method’s name denote the training loss. Models with the -d- middle part estimate full diagonal of the decoder variance, -s- estimate only a scalar, and -c- use a fixed scalar variance as a hyperparameter. All variants are using the standard Gaussian latent model. Models using the VampPrior are denoted by extending the decoder variance symbol by the letter v-, i.e. -dv-, -sv-, -cv-. The last part of the name denotes score, -rs stands for the sampled rec. probability (7) with  $L = 100$ , -rm for (8), -el for the ELBO (3) composed of -rs and KLD, -jc for (10), -di for (14).

the x-axis will vary. In the averaged ranks, VAE is on the Pareto front in both fit and prediction times, see Fig. 5b. Its prediction complexity is given mainly by the choice of the number of samples taken in the computation of the sampled reconstruction score (7). The kNN detector has negligible training time, given only by the construction of the tree structure representing data, but seems to be mostly unusable on image data due to slow prediction times on datasets that are large in dimension and number of samples. The fmGAN finds itself in a completely reversed scenario.

#### D. Other influences

In this section, we report results of three sources of variability of performance of AD methods that were found to have minimal impact.

*Ensembles/Hyperparameter averaging* The benefits of ensembles in prior art seem to be mixed. While [17] claims that a

combination of VAE or GAN ensembles using WAIC might be useful, Ref. [59] claims a negligible effect. In our experiments, we have used ensembles as a way to reduce uncertainty in hyperparameters [90], meaning that unlike in [17], models in ensembles were of a single type differing only in architecture. The effect of such an ensemble on average AUC was overall zero, sometimes even negative, see Supplementary E-C, Table E.4. Exceptions are methods based on GANs featuring improvements by 0.02 in average AUC. These findings are on par with those in [59].

*Bayesian optimization* Bayesian optimization was introduced as an alternative to random search of hyperparameters. It has the potential to find better hyperparameters with a low number of trials. Comparison of the random search and Bayesian optimization under the same conditions is reported in the Supplementary, Table E.3. We can conclude that Bayesian optimization was able to find hyperparameters with better performance for the vast majority of the methods. However, this improvement is often quite small and insufficient to change the ranks of the methods. Notable exceptions are the GANomaly and GAN methods which improved by two ranks.

*Performance metric AUC/TPR* The hypothesis that the methods may perform differently when chosen for different optimization criteria than the usual AUC has not been proved. The results are summarized in the Supplementary, Tab. E.1. While small changes in performance have been observed, the ranks of the methods remain the same in both criteria, AUC and TPR@5%.

## VI. CONCLUSION

The presented extensive comparison of anomaly detection methods based on deep generative methods, namely variants of flows, variational autoencoders and generative adversarial networks with methods based on alternative paradigms (Support vector machines, random forests, histogram and distance-based methods) revealed that the performance of anomaly detection methods strongly depends on experimental conditions. We have identified the main sources of variation to be the choice of data and hyperparameter tuning. We presented detailed results under a combination of experimental conditions, called contexts, specifically the data context, hyperparameter context and economic context.

In the dataset context, a clear distinction in performance was found between tabular data, image data containing statistical and semantic anomalies. The main distinction in hyperparameter selection was how many anomalies were present in the validation set. Different order of performance of the tested method was observed for a different combination of dataset type and hyperparameter selection context, explaining various outcomes of comparison in the prior art. We strongly recommend to authors to provide details on the context of their experimental studies in the future.

There are many aspects that have not been covered and remain a topic of future works, such as designing ensembles of methods of various types, or identification of the relevant kind of anomalies in the semantic datasets.

The comparison is not only aimed at researchers, but also at practitioners desiring accurate methods with low computa-



- bert, Ira Assent, and Michael 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.
- [13] Raghavendra Chalapathy and Sanjay Chawla. Deep learning for anomaly detection: A survey. *arXiv:1901.03407 [cs]*, 2019.
- [14] Raghavendra Chalapathy, Aditya Krishna Menon, and Sanjay Chawla. Anomaly detection using one-class neural networks. *arXiv:1802.06360 [cs]*, 2019.
- [15] Raghavendra Chalapathy, Edward Toth, and Sanjay Chawla. Group Anomaly Detection using Deep Generative Models. *arXiv:1804.04876 [cs]*, April 2018.
- [16] Xiaoran Chen and Ender Konukoglu. Unsupervised detection of lesions in brain MRI using constrained adversarial auto-encoders. *arXiv:1806.04972 [cs]*, 2018.
- [17] Hyunsun Choi, Eric Jang, and Alexander A. Alemi. WAIC, but Why? Generative Ensembles for Robust Anomaly Detection. *arXiv:1810.01392 [cs, stat]*, May 2019.
- [18] Bin Dai and David Wipf. Diagnosing and enhancing vae models. *arXiv:1903.05789 [cs]*, 2019.
- [19] Lucas Deecke, Robert Vandermeulen, Lukas Ruff, Stephan Mandt, and Marius Kloft. Image anomaly detection with generative adversarial networks. In *Joint european conference on machine learning and knowledge discovery in databases*, pages 3–17. Springer, 2018.
- [20] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*, 7(Jan):1–30, 2006.
- [21] Madson L. D. Dias, César Lincoln C. Mattos, Ticiana L. C. da Silva, José Antônio F. de Macedo, and Wellington C. P. Silva. Anomaly Detection in Trajectory Data with Normalizing Flows. *arXiv:2004.05958 [cs, stat]*, April 2020.
- [22] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. *arXiv:1605.08803 [cs]*, 2017.
- [23] Lori E Dodd and Margaret S Pepe. Partial auc estimation and regression. *Biometrics*, 59(3):614–623, 2003.
- [24] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv:1605.09782 [cs]*, 2017.
- [25] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [26] Andrew F. Emmott, Shubhomoy Das, Thomas Dietterich, Alan Fern, and Weng-Keen Wong. Systematic construction of anomaly detection benchmarks from real data. In *Proceedings of the ACM SIGKDD workshop on outlier detection and description*, pages 16–21, 2013.
- [27] Tolga Ergen and Suleyman Serdar Kozat. Unsupervised anomaly detection with LSTM neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 31(8):3127–3141, Aug 2020.
- [28] Gilberto Fernandes, Joel JPC Rodrigues, Luiz Fernando Carvalho, Jalal F Al-Muhtadi, and Mario Lemes Proença. A comprehensive survey on network anomaly detection. *Telecommunication Systems*, 70(3):447–489, 2019.
- [29] Manuel Fernández-Delgado, Eva Cernadas, Senén Barro, and Dinani Amorim. Do we need hundreds of classifiers to solve real world classification problems? *The journal of machine learning research*, 15(1):3133–3181, 2014.
- [30] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: Masked Autoencoder for Distribution Estimation. *arXiv:1502.03509 [cs, stat]*, June 2015.
- [31] Markus Goldstein and Andreas Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly detection algorithm. *KI-2012: Poster and Demo Track*, pages 59–63, 2012.
- [32] Markus Goldstein and Seiichi Uchida. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PloS one*, 11(4), 2016.
- [33] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [34] Parikshit Gopalan, Vatsal Sharan, and Udi Wieder. PID-Forest: Anomaly Detection via Partial Identification. *arXiv:1912.03582 [cs]*, 2019.
- [35] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. FFJORD: Free-form Continuous Dynamics for Scalable Reversible Generative Models. *arXiv:1810.01367 [cs, stat]*, October 2018.
- [36] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017.
- [37] Stefan Harmeling, Guido Dornhege, David Tax, Frank Meinecke, and Klaus-Robert Müller. From outliers to prototypes: ordering data. *Neurocomputing*, 69(13-15):1608–1618, 2006.
- [38] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [39] Tim Head, Manoj Kumar, Holger Nahrstaedt, Gilles Louppe, and Iaroslav Shcherbatyi. scikit-optimize, 2018. 10.5281/zenodo.1157319.
- [40] Yongjun Hong, Uiwon Hwang, Jaeyoon Yoo, and Sungroh Yoon. How generative adversarial networks and their variants work: An overview. *ACM Computing Surveys (CSUR)*, 52(1):1–43, 2019.
- [41] Mike Innes. Flux: Elegant machine learning with julia. *Journal of Open Source Software*, 2018.
- [42] Tomoharu Iwata and Yuki Yamanaka. Supervised Anomaly Detection based on Deep Autoregressive Density Estimators. *arXiv:1904.06034 [cs, stat]*, April 2019.
- [43] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980 [cs]*, 2017.
- [44] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative Flow with Invertible 1x1 Convolutions. *arXiv:1807.03039 [cs, stat]*, July 2018.
- [45] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv:1312.6114 [stat]*, 2014.

- [46] B Ravi Kiran, Dilip Mathew Thomas, and Ranjith Parakkal. An overview of deep learning based methods for unsupervised and semi-supervised anomaly detection in videos. *Journal of Imaging*, 4(2):36, 2018.
- [47] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Why normalizing flows fail to detect out-of-distribution data. *arXiv:2006.08545 [stat]*, 2020.
- [48] Mark Kliger and Shachar Fleishman. Novelty detection with gan. *arXiv:1802.10560 [cs]*, 2018.
- [49] Ivan Kobyzev, Simon J. D. Prince, and Marcus A. Brubaker. Normalizing Flows: An Introduction and Review of Current Methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020.
- [50] Hans-Peter Kriegel, Matthias Schubert, and Arthur Zimek. Angle-based outlier detection in high-dimensional data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 444–452, 2008.
- [51] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [52] Donghwoon Kwon, Hyunjoon Kim, Jinho Kim, Sang C Suh, Ikkyun Kim, and Kuinam J Kim. A survey of deep learning-based network anomaly detection. *Cluster Computing*, pages 1–13, 2019.
- [53] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- [54] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 Eighth IEEE International Conference on Data Mining*, pages 413–422. IEEE, 2008.
- [55] Yezheng Liu, Zhe Li, Chong Zhou, Yuanchun Jiang, Jianshan Sun, Meng Wang, and Xiangnan He. Generative adversarial active learning for unsupervised outlier detection. *IEEE Transactions on Knowledge and Data Engineering*, 2019.
- [56] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv:1511.05644 [cs]*, 2016.
- [57] Nour Moustafa, Jiankun Hu, and Jill Slay. A holistic review of network anomaly detection systems: A comprehensive survey. *Journal of Network and Computer Applications*, 128:33–55, 2019.
- [58] Norman Mu and Justin Gilmer. MNIST-C: A Robustness Benchmark for Computer Vision. *arXiv:1906.02337 [cs]*, June 2019.
- [59] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do Deep Generative Models Know What They Don’t Know? *arXiv:1810.09136 [cs, stat]*, February 2019.
- [60] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisaccho, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
- [61] Guansong Pang, Longbing Cao, Ling Chen, and Huan Liu. Learning Representations of Ultrahigh-dimensional Data for Random Distance-based Outlier Detection. *arXiv:1806.04808 [cs, stat]*, June 2018.
- [62] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton van den Hengel. Deep learning for anomaly detection: A review. *arXiv:2007.02500 [cs]*, 2020.
- [63] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing Flows for Probabilistic Modeling and Inference. *arXiv:1912.02762 [cs, stat]*, December 2019.
- [64] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked Autoregressive Flow for Density Estimation. *arXiv:1705.07057 [cs, stat]*, June 2018.
- [65] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [66] Joao Pereira and Margarida Silveira. Unsupervised anomaly detection in energy time series data using variational recurrent autoencoders with attention. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 1275–1282. IEEE, 2018.
- [67] Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. Ocgan: One-class novelty detection using GANs with constrained latent representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2898–2906, 2019.
- [68] Tomáš Pevný. Loda: Lightweight on-line detector of anomalies. *Machine Learning*, 102(2):275–304, 2016.
- [69] Tomas Pevny, Vasek Smidl, Martin Trapp, Ondrej Polacek, and Tomas Oberhuber. Sum-product-transform networks: Exploiting symmetries using invertible transformations. *arXiv:2005.01297 [stat]*, 2020.
- [70] Stanislav Pidhorskyi, Ranya Almohsen, and Gianfranco Doretto. Generative probabilistic novelty detection with adversarial autoencoders. In *Advances in neural information processing systems*, pages 6822–6833, 2018.
- [71] Marco AF Pimentel, David A Clifton, Lei Clifton, and Lionel Tarassenko. A review of novelty detection. *Signal Processing*, 99:215–249, 2014.
- [72] Emanuele Principi, Fabio Vesperini, Stefano Squartini, and Francesco Piazza. Acoustic novelty detection with adversarial autoencoders. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 3324–3330. IEEE, 2017.
- [73] Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. Efficient algorithms for mining outliers from large data sets. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 427–438, 2000.
- [74] Jie Ren, Peter J. Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark A. DePristo, Joshua V. Dillon, and Balaji Lakshminarayanan. Likelihood Ratios for Out-of-Distribution Detection. *arXiv:1906.02845 [cs, stat]*, December 2019.
- [75] Lukas Ruff, Jacob R Kauffmann, Robert A Vander-

- meulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *arXiv:2009.11732 [cs]*, 2020.
- [76] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402, 2018.
- [77] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *Advances in neural information processing systems*, pages 2234–2242, 2016.
- [78] Thomas Schlegl, Philipp Seeböck, Sebastian M. Waldstein, Georg Langs, and Ursula Schmidt-Erfurth. F-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks. *Medical Image Analysis*, 54:30–44, May 2019.
- [79] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International Conference on Information Processing in Medical Imaging*, pages 146–157. Springer, 2017.
- [80] Maximilian Schmidt and Marko Simic. Normalizing flows for novelty detection in industrial time series data. *arXiv:1906.06904 [cs, stat]*, June 2019.
- [81] Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7):1443–1471, 2001.
- [82] Vít Škvára, Tomáš Pevný, and Václav Šmídl. Are generative deep models for novelty detection truly better? *arXiv:1807.05027 [cs]*, 2018.
- [83] Vít Škvára, Václav Šmídl, Tomáš Pevný, Jakub Seidl, Aleš Havránek, and David Tskhakaya. Detection of alfvén eigenmodes on COMPASS with generative neural networks. *Fusion Science and Technology*, 2020.
- [84] Václav Šmídl, Jan Bím, and Tomáš Pevný. Anomaly scores for generative models. *arXiv:1905.11890 [stat]*, 2019.
- [85] Ingo Steinwart, Don Hush, and Clint Scovel. A classification framework for anomaly detection. *Journal of Machine Learning Research*, 6(8):211–232, 2005.
- [86] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. *arXiv:1711.01558 [stat]*, 2019.
- [87] Jakub Tomczak and Max Welling. VAE with a Vamp-Prior. In *International Conference on Artificial Intelligence and Statistics*, pages 1214–1223, 2018.
- [88] Hongzhi Wang, Mohamed Jaward Bah, and Mohamed Hammad. Progress in outlier detection techniques: A survey. *IEEE Access*, 7:107964–108000, 2019.
- [89] Xuhong Wang, Ying Du, Shijie Lin, Ping Cui, Yuntian Shen, and Yupu Yang. advae: A self-adversarial variational autoencoder with gaussian anomaly prior knowledge for anomaly detection. *Knowledge-Based Systems*, 190:105187, 2020.
- [90] Andrew Gordon Wilson. The case for bayesian deep learning. *arXiv:2001.10995 [cs]*, 2020.
- [91] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv:1708.07747 [cs]*, 2017.
- [92] Haowen Xu, Wenxiao Chen, Nengwen Zhao, Zeyan Li, Jiahao Bu, Zhihan Li, Ying Liu, Youjian Zhao, Dan Pei, Yang Feng, et al. Unsupervised anomaly detection via variational auto-encoder for seasonal kpis in web applications. In *Proceedings of the 2018 World Wide Web Conference*, pages 187–196, 2018.
- [93] Masataka Yamaguchi, Yuma Koizumi, and Noboru Harada. Adaflow: Domain-adaptive density estimator with application to anomaly detection and unpaired cross-domain translation. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3647–3651. IEEE, 2019.
- [94] Rong Yao, Chongdang Liu, Linxuan Zhang, and Peng Peng. Unsupervised Anomaly Detection Using Variational Auto-Encoder based Feature Extraction. In *2019 IEEE International Conference on Prognostics and Health Management (ICPHM)*, pages 1–7, San Francisco, CA, USA, June 2019. IEEE.
- [95] Houssam Zenati, Chuan Sheng Foo, Bruno Lecouat, Gaurav Manek, and Vijay Ramaseshan Chandrasekhar. Efficient GAN-based anomaly detection. *arXiv:1802.06222 [cs]*, 2019.
- [96] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Information maximizing variational autoencoders. *arXiv:1706.02262 [cs]*, 2018.
- [97] Yue Zhao, Zain Nasrullah, and Zheng Li. Pyod: A python toolbox for scalable outlier detection. *Journal of Machine Learning Research*, 20(96):1–7, 2019.
- [98] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International Conference on Learning Representations*, 2018.

# Supplementary material for: Comparison of Anomaly Detectors: Context Matters

## APPENDIX A EXPERIMENTAL SETUP

### A. Data splits

Our decision for not doing repeated experiments on different random crossvalidation data splits in MNIST, SVHN2, FashionMNIST and CIFAR10 datasets was motivated by a need to reduce the computational requirements of experiments and justified by preliminary results that are summarized in Tab. A.1. The consistent performance across different resamplings of the datasets is due to their large size and homogeneity.

On the other hand in the Tab. A.2 we can see that in the case of tabular data and the other image datasets the performance on each split varies to a greater extent. The biggest changes can be observed on MVTec-AD datasets, due to its small size, whereas MNIST-C shows only small changes. On tabular data the situation is even more chaotic due to the dataset diversity, thus using only one seed for the random split is in this case unreliable.

### B. Hyperparameter settings

An overview of sampled hyperparameters for all models used in our experiments is in Tab. A.3–A.5. The default hyperparameters of classical methods used in the context of validation set without anomalies are present in Tab. A.6

## APPENDIX B DATASET CONTEXT DETAILS

The underlying data for critical diagrams in Section V-A are displayed in Tables B.1, B.2, B.3, and B.4 in the form of testing AUC for all individual datasets and all tested methods.

## APPENDIX C EXTENDING IMAGE DATASETS RESULTS

In Sec. V-A we stated that some models perform poorly on datasets with semantic anomalies, in which there are multiple sources of variation and the true anomalous information is pronounced in few of them. It is often the case that the methods focus more on the other sources of variation mainly the statistical aspects like background. As an entertaining example we have devised a simple anomaly detector *blpix* for detecting the airplane class of the CIFAR10 dataset. As the name suggest it sums up all blue pixels in an image, resulting in anomaly score  $s_{\text{blpix}}(x) = -\sum_{i,j} 1_{x_{i,j}=\text{blue}}$ . Despite its simplicity, it achieves 0.68 AUC on both the validation and the test split, which is comparable with other models, see Fig. C.1.

More visual examples of anomalies for the best performing models on MNIST, FashionMNIST, CIFAR10 and SVHN2 datasets are shown in Fig. C.2–C.3. Similar examples are also shown for MVTec-AD and MNIST-C datasets in figures C.4 and C.5

## APPENDIX D HYPERPARAMETER CONTEXT

Following the extension of image results, here we provide insights into the sensitivity of methods to the number of anomalies as seen from the point of view of different image datasets, which further illustrates the dataset context as an important factor in the model ranking. In Fig. D.1 we show the analogy of Fig. 3 for a more granular grouping of image datasets. The first thing that stands out that the performance on MVTec-AD dataset is constant while increasing the number of labeled samples from 0.01% to 0.1%. This is due to the size of the dataset, which at these percentages does not contain enough samples to choose the best hyperparameters of each method and the choice is essentially arbitrary but consistent within this range. We admit that such points should not be even reported in the figures, but we kept it there for consistency with other datasets that have more samples. Although this behaviour is present also in the case of the tabular data, where the sample size are as low as  $\sim 80$  with Breast Tissue, the presence of large datasets such as the KDD99(10%) with  $\sim 500k$  samples allows us to see the general trend more clearly.

The second thing that we have already touched upon many times is the fact that with each dataset group the best performing models at both ends are different, with the exception of fmGAN, which does not work in the case of clean validation context (left-most points).

Lastly we would like to shed some light on the use of the PR@% metric in the context of the plots in question. Though throughout the text we have used AUC metric exclusively, the problem with using AUC on a portion of most anomalous samples is that we have a high chance of encountering samples of one class only, where the metric is not well defined. Therefore we opted for precision, where we assume that the threshold for discrete prediction is always lower than the last sample, as sorted by the anomaly score and thus all of them are labeled as positives.

*Sensitivity study of OC-SVM:* In order to better support the thesis about insufficient tuning of the OCSVM method, we provide here the average ranks of all methods trained on tabular data, where the fully optimized method was replaced by its OCSVM-RBF counterpart (only the width of the rbf kernel was sampled and hyperparameter  $\nu$  was set to 0.5). As average ranks are largely dependent also on the composition of the methods this change has a profound effect on ranks of other methods such as VAE, which takes the lead, see Tab. D.1.



	seed = 1	seed = 2	seed = 3	seed = 4	seed = 5
cifar10	<b>0.63</b>	<b>0.63</b>	<b>0.63</b>	0.62	<b>0.63</b>
fmnist	0.77	0.77	0.77	<b>0.79</b>	0.77
mnist	<b>0.88</b>	<b>0.88</b>	0.87	0.87	<b>0.88</b>
svhn2	<b>0.54</b>	<b>0.54</b>	0.53	<b>0.54</b>	<b>0.54</b>

(a) VAE

	seed = 1	seed = 2	seed = 3	seed = 4	seed = 5
cifar10	<b>0.65</b>	<b>0.65</b>	0.64	<b>0.65</b>	<b>0.65</b>
fmnist	<b>0.82</b>	<b>0.82</b>	<b>0.82</b>	<b>0.82</b>	0.81
mnist	0.88	<b>0.89</b>	0.88	0.88	0.88
svhn2	<b>0.52</b>	<b>0.52</b>	<b>0.52</b>	<b>0.52</b>	<b>0.52</b>

(b) GANomaly

TABLE A.1: Summary tables for the VAE and GANomaly models on different seeds of the image datasets. Reported values are the AUC metrics on the test data, averaged over 10 anomaly classes.

	aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	osvm	pidf	maf	rnvp	sptn	fmgm	gan	mgal	dagm	dsvd	rpn	vaek	vaeo
	6.9	13.4	10.3	6.5	7.0	11.4	14.9	14.5	8.4	16.1	11.4	<b>2.9</b>	14.3	8.9	9.0	10.7	11.4	12.0	22.3	19.8	16.2	12.1	11.0	6.8
SEED 1	9.1	13.6	10.8	6.6	8.0	12.7	15.0	12.6	9.4	14.4	11.5	<b>3.8</b>	13.6	8.6	9.5	11.7	10.0	11.2	18.2	20.4	17.3	11.6	11.8	6.4
SEED 2	7.4	14.2	10.0	6.0	7.5	12.7	16.4	14.2	10.4	15.9	11.8	<b>5.0</b>	14.5	8.3	7.7	10.8	11.2	10.8	18.4	20.2	17.0	11.3	10.9	6.0
SEED 3	6.9	14.5	10.5	6.3	7.0	12.8	15.8	14.4	9.4	14.6	12.3	<b>3.9</b>	14.1	8.5	8.6	11.9	11.0	10.8	16.6	19.7	17.1	12.1	11.4	5.6
SEED 4	8.5	12.0	8.4	7.3	6.6	13.2	16.2	13.9	9.9	15.9	12.4	<b>4.7</b>	13.7	8.0	7.8	12.6	11.2	10.4	19.0	19.4	16.5	10.4	11.0	6.6
SEED 5	6.8	14.0	9.4	6.8	7.8	12.9	15.4	13.8	8.8	15.2	12.4	<b>3.9</b>	14.3	9.1	9.4	12.6	11.6	10.8	17.6	20.1	16.8	11.5	10.6	6.0

(a) tabular

	aae	gano	skip	vae	wae	knn	osvm	fano	fmgm	dsvd	vaek	vaeo
	4.1	9.9	10.6	2.9	<b>2.4</b>	5.4	6.0	3.2	3.1	6.6	4.5	<b>2.4</b>
SEED 1	3.9	9.9	10.6	2.9	<b>2.1</b>	5.7	5.9	3.3	2.8	6.6	4.7	2.3
SEED 2	3.9	9.9	10.6	3.0	<b>2.2</b>	5.4	6.0	3.0	2.7	6.6	4.6	2.8
SEED 3	3.7	9.9	10.7	3.1	<b>2.0</b>	5.7	6.0	3.1	2.4	6.8	4.6	2.6
SEED 4	3.7	9.9	10.6	2.9	2.3	5.4	6.1	3.0	<b>2.1</b>	6.2	4.5	2.3
SEED 5	3.9	9.9	10.6	2.9	<b>2.1</b>	5.5	6.0	3.0	3.0	6.0	4.6	2.6

	aae	gano	skip	vae	wae	knn	osvm	fano	fmgm	dsvd	vaek	vaeo
	4.7	7.0	9.0	<b>3.0</b>	3.7	6.0	6.7	4.0	4.3	9.7	9.0	8.0
SEED 1	4.0	9.0	8.7	<b>2.3</b>	3.7	4.7	7.3	5.0	6.3	9.0	8.7	7.0
SEED 2	2.7	8.3	8.3	2.7	6.3	7.0	7.0	<b>1.3</b>	3.0	8.0	9.7	10.3
SEED 3	4.0	9.3	6.3	<b>2.0</b>	4.0	6.0	7.7	3.3	5.3	6.7	9.7	9.7
SEED 4	3.7	9.7	8.7	4.3	3.0	6.3	5.3	6.7	<b>2.7</b>	7.0	9.0	9.3
SEED 5	3.7	9.0	8.3	2.3	<b>1.3</b>	5.7	7.7	4.3	4.7	7.7	9.3	11.0

(b) MNIST-C

(c) MVTec-AD

TABLE A.2: Comparison of average ranks in AUC on different train/val/test splits of datasets that were subject to repeated experiments. The first row corresponds to the average rank in AUC, which has been averaged over all 5 repetitions.

model	parameter	value set
ABOD	$n$	$\{1, 2, \dots, 100\}$
	method	fast
HBOS	no. bins	$\{2, 4, \dots, 100\}$
	$\alpha$	$\{0.05, 0.1, \dots, 1\}$
	tolerance	$\{0, 0.05, \dots, 1\}$
IF	no. estimators	$\{50, 100, \dots, 500\}$
	max samples	$\{0.5, 0.6, \dots, 1\}$
	max features	$\{0.5, 0.6, \dots, 1\}$
KNN	$n$	$\{1, 3, \dots, 101\}$
LODA	no. bins	$\{2, 4, \dots, 100\}$
	no. cuts	$\{40, 60, \dots, 500\}$
LOF	$n$	$\{1, 2, \dots, 100\}$
OC-SVM	$\gamma$	$10^x, x \in \{-4, -3.9, \dots, 2\}$
	$\nu$	$\{0.01, 0.5, 0.99\}$
	kernel	$\{\text{rbf, sigmoid, polynomial}\}$
PIDForest	max depth	$\{6, 8, \dots, 10\}$
	no. trees	$\{50, 75, \dots, 200\}$
	max samples	$\{50, 100, 250, 500, 1000, 5000\}$
	max buckets	$\{3, 4, 5, 6\}$
	$\epsilon$	$\{0.05, 0.1, 0.2\}$

TABLE A.3: Hyperparameters of the classical models.

## APPENDIX E

### OTHER INFLUENCES

#### A. Performance metrics

Throughout the main text, we have reported only results for the AUC metric. Though this has had quite practical

implications on the degrees of freedom, we are aware that the choice of metric may favour each method differently. To this extend we provide Tab. E.1 containing simple comparison of ranks as measured in AUC and TPR@5% of FPR on each of the dataset types. The changes in ranks are more pronounced in the case of tabular datasets, where the biggest improvement when going from AUC to TPR@5% is seen with the worst performing model, MOGAAL, thus swapping its place with DAGMM. The biggest fall has been experienced by the WAE, VAE method but not as much to affect the overall ranking at the top positions. On image datasets the changes are really minor.

#### B. Bayesian optimization

Bayesian optimization allowed us to fill the gaps in the discrete hyperparameter sampling ranges described in tables A.3–A.5, though we have not always included the range as a whole, but for example in case of the neural network architectures we used powers of 2 and optimized the exponent. As we have been building on top of already trained models, we have not found an easy way how to optimize hyperparameters such as weight decay, where we used to indicate 0.0 as not using this regularization, while still using log-uniform with base 10 for sampling, therefore we have chosen to keep the regularization while enlarging the parameter space to smaller values. In spite of this the definition of sampling ranges was generally a simple exercise and the biggest hurdle was on the implementation side

of things, which we will discuss briefly in the implementation details section F-F.

Regarding the hyperparameters for the optimization itself, we used the default parameters as provided by the scikit-optimise framework, see Tab. E.2, which is certainly orthogonal to one of the conclusions of the main text, however including another degree of freedom to this large scale comparison was computationally unrealistic. Due to the same constraints, we have limited the Bayesian optimization to tabular data only, as the biggest time constraint is given by the sequential nature of training, whereas random sampling allows to fit multiple hyperparameters on the same dataset in parallel. kNN, ABOD, LOF methods have been purposely left out of this experiment as their hyperparameter space is one dimensional and thus extensively sampled already. In the case of the MOGAAL methods the decision has been made based on the nature of available hyperparameters, which affected only the training part but not the model’s architecture, therefore we did not expect much improvement there.

As has been already mentioned we have trained initially 50 random samples (with already precomputed methods we have just selected the random search runs) and then used the Bayesian optimization to propose the next points, while updating the underlying process. At the end we have been left with 100 samples, which we have evaluated in the same way as in the case of random sampling, i.e. choosing the best performing models based on AUC on validation data averaged over 5 repetitions, which also served as the objective function optimized for the Bayesian optimization. Comparison of ranks between random sampling and Bayesian optimization is provided in Tab. E.3. We can see that generally all methods improved in terms of average change in AUC, with the exception of adVAE, whose performance degraded on 21 out of 40 datasets. On the other hand the methods that did improve significantly was GANomaly, which jumped by 3 ranks on average, surpassing AAE and flow models just to name a few. Other notable improvement has been made by GAN and RealNVP flow.

Some of the drops in performance may be explained by our choice of default parameters, as for some model–dataset combination the optimization did not explore as much of the hyperparameter space or even got stuck at one point. As a result, the pool of methods may have shrunk close to the original 50 random samples. One other factor that may have affected the results, has been the lack of hyperparameters restrictions that we have manually imposed after random sample has been drawn, such having latent dimension of autoencoders lower than the input dimension.

When put together, the individual gains we have seen (namely GANomaly, flows and gans) are definitely worth exploring, when one has only one objective to aim at and enough computational budget to explore even the hyperparameters of Bayesian optimization itself.

### C. Ensembles experiments

Creation of ensembles/hyperparameter averaging method consisted of two steps. Firstly the performance metrics were

precomputed for all the methods and each repetition or anomaly class, in order to quickly assert which hyperparameters to include. Secondly we created multiple ensembles for each combination of the criterion metric and size, totaling in 4 different ensembles per each repetition, which allowed us to compare individual effects. Resulting average ranks in the AUC metric and average improvement in the AUC are presented in odd/even rows respectively of tables E.4.a–c for each dataset type and combination of criterion and ensemble size. We can see that in the majority the effect on AUC is zero, with notable exception of GANomaly on tabular data and feature matching GAN on images containing semantic anomalies. Notable observation is that smaller ensembles perform usually better, which may suggest that by not weighting individual members of the ensembles or by not employing some heuristic for choosing the optimal size, we are crudely including hyperparameters, that only hinder the performance.

## APPENDIX F IMPLEMENTATION DETAILS

### A. Classical methods

Due to the quadratic scaling of OC-SVM in the number of training samples, we have implemented the OC-SVM on image data as an ensemble of 10 OC-SVM models with the same hyperparameters but trained on 10 equal-sized subsets of training data. The resulting anomaly score of such an ensemble was obtained by averaging.

For the kNN model, we have used 3 anomaly scores that differed at the distance computation, as described in [37]. As opposed to the other methods, implementation of kNN required further downscaling of MVTEC-AD images to  $64 \times 64 \times 3$  in order to fit the computational envelope.

The original PIDForest [34] method’s implementation provided scikit-learn model API, however that’s where the convenience ended. We have encountered a number of errors both during training and prediction, some of which we have identified were caused by features with constant values, which we had to filter out from the splits, before proceeding. It’s quite unfortunate when, what seems to be a sound method on paper, is handicapped by a subpar implementation, which is both slow and breaks at the first sight of the unknown. Though the authors have promised a better one, as of the time of writing, such an option was not available. As opposed to the original implementation we used negative sparsity as the anomaly score and parametrized its quantiles at three levels  $\{0.10, 0.25, 0.50\}$ .

### B. Autoencoders & GANs

The training of VAE, WAE and AAE models has been stopped prematurely if the loss on normal samples in the validation dataset has not improved for 200 batches. For GAN and fmGAN, patience of 50 batches was used and the discriminator loss was used for early stopping. We have not used batch normalization in any of these models, as it is recommended not to do so due to the additional stochasticity that it introduces in the training process.

In the construction of VAE, WAE, AAE, GAN and fmGAN models on image data, convolution layers are used. First the number of convolution layers  $n_c$  was sampled. To construct decoder/generator parts, the first  $n_c$  elements of the channels, kernelsizes and scalings vectors in Tab. A.4 were used, e.g. for  $n_c = 2$ , a decoder with (16, 32) channels, (3, 5) kernelsizes and (1, 2) scaling was constructed, while the architecture of the encoder/discriminator was flipped. Downscaling in encoders/discriminators was done with maxpooling, upscaling in decoders/discriminators was done via transposed convolutions.

Each of the  $k$  components of a VampPrior initialized with an average sample of the training data, perturbed with noise sampled from  $\mathcal{N}(0, 1)$  in each dimension.

The most challenging implementation was that of ad-VAE [89] because no training loop or loss functions was provided in the code provided by the author. Additionally, the other implementation recommended by the author differs from the actual paper in terms of loss functions and some hyperparameters. In spite of these uncertainties, we have managed to implement a working version of the model that follows the paper as closely as possible. The only difference is the architecture of the encoder and decoder.

In the case of GANomaly and skip-GANomaly models for image data, we implemented the exact same architectures as in original papers [3, 4] and the only addition to those models is weight decay. This architecture, which is pretty much similar for both models, has one small disadvantage. The width of the input image must be equal to its height and that must be divisible by 32. This allows fully convolution encoders to compress input image into 1D latent vector resp.  $(1, 1, \dim(z))$  tensor. Therefore we have to resize MNIST and FashionMNIST datasets in the preprocess phase from shape 28x28 to 32x32.

To avoid overfitting of deep models we decided to use an early stopping procedure. The criterion for this procedure most often corresponds to the anomaly score formula or loss function, so it is different for every model. For example GANomaly’s early stopping criterion is simple generator loss  $L_G = w_1 L_{adversarial} + w_2 L_{contextual} + w_3 L_{latent}$  [3], however for skip-GANomaly we use weighted average of generator loss  $L_G$  and discriminator loss  $L_D$  because both those terms are also part of anomaly score. Weights in  $L_G$  are typically treated as hyperparameters which we sample. Exception is image version of GANomaly where we used fixed weights  $w_1 = 1$ ,  $w_2 = 50$  and  $w_3 = 1$  since those weights were empirically derived in original paper as optimal for the same datasets we use.

During the implementation of fAnoGAN we were forced to make some adjustments to the original code. Firstly we switched Wasserstein GAN with gradient penalty (WGAN-GP [36]) to its predecessor WGAN with clipping weights [6] because we were not able to compute the second order derivative needed to update weights of the model due to Julia limitations. Later we tried Pytorch implementation of WGAN-GP, but it performed on par with WGAN so we dropped it from results.

Secondly, we replaced the original ResNet [38] generator and discriminator with much lighter architecture (in the same

way as with GANomaly), that we also use for autoencoders. This helps to compare models more fairly, as their performance now depends mainly on the model concept, not the architecture itself. In addition, fAnoGAN is the only model in this category without an early stopping procedure, because its loss behavior is not standard and we had a hard time figuring out when to stop it.

### C. Normalizing flows

While implementing RealNVP and MAF flows on tabular data we tried to follow as closely as possible the code that has accompanied [64], however on a closer inspection, we found that there is a lot of variety in the literature, which we tried to include into hyperparameters of each of the methods. This variety comes mainly from the different techniques, that authors [22, 35, 44] employ to battle instabilities in the training of deep flows, though the exact nature of which is often not disclosed. In our case, the root of most problems was in scale conditioner’s output saturating (yielding Inf) in the subsequent exponentiation. We have been able to isolate this to specific samples in the training datasets such as Statlog Shuttle, which unfortunately contained outliers with one of the features having disproportional values. The MAF flows seemed to be more affected by this due to their autoregressive structure, which propagated the discrepancy into other features in each subsequent flow.

Using batch normalization as suggested by [64] should supposedly improve both stability and target likelihoods, which we could confirm while experimenting on small problems, however it did not help us with the aforementioned issues. We’ve found that batch normalization with floating averages as used in [22] performed poorly even on smaller problems, thus we used only statistics computed for each batch individually. While looking deeper into the performance of different hyperparameters, we have found that none of the models on KDD99 (10%) with batch normalization trained properly yielding AUC’s below 0.5. As per the original implementation, during evaluation the statistics in batch normalization were initialized from the whole training dataset.

As suggested by [22] we have also experimented with initializing the flow to identity by zeroing the last layers of conditioner networks, however it also did not help. Furthermore in case of RealNVP flow we have tried to use ”tanh normalization” together with learn-able location  $\alpha_s \in \mathbb{R}$  and scale  $\beta_s \in \mathbb{R}$  parameters

$$\alpha_s \oplus \exp \beta_s \odot \tanh (s), \quad (15)$$

where  $s$  is the scale conditioner’s output and  $\oplus$ ,  $\odot$  are element-wise operations. Using this scaling together with tanh activations, we were able to train even deeper and wider architectures.

One of the degrees of freedom is the choice of conditioner architecture. Where on images the consensus [22, 44] seems to be to use one convolution neural network for both location and scale parameters and feeding the network positive and negative input. In the case of tabular flows, where the number of weights is not as closely controlled variable, the additional

options are to employ either two networks for both scale and location (RealNVP implementation in [64]), or one network whose output is fed through another pair of single-layer networks to provide the two outputs (MAF implementation in [64]). We have opted for using two separate networks for both RealNVP/MAF flows to equal the situation.

Even the exact formulation transformation function  $f$  from (1) has different forms in the chosen class of autoregressive affine flows [63], across different implementations. We have used the following definition of forward step  $f^{-1}$  for RealNVP

$$\begin{aligned} z_{\leq d} &= \mathbf{x}_{\leq d}, \\ z_{> d} &= \exp\left(-\frac{1}{2}\mathbf{s}\right) \odot (\mathbf{t} - \mathbf{x}_{\leq d}), \end{aligned} \quad (16)$$

where  $\mathbf{s}$ ,  $\mathbf{t}$  are outputs of scale and location conditioners respectively when given the first half of the input  $\mathbf{x}_{\leq d} = \mathbf{z}_{\leq d}$ . In case of MAF flow the forward step had the following form

$$z_i = \exp\left(-\frac{1}{2}s_i\right) (t_i - x_i), \quad \forall i \in \{1, \dots, D\} \quad (17)$$

where  $s_i$ ,  $t_i$  are the  $i$ -th elements of outputs of the autoregressive conditioners (MADE networks [30]) for scale and location parameters respectively, when given the whole input vector  $\mathbf{x}$ . It is important to stress the halving of the outputs of the scale conditioners, as this is another means with which to battle the aforementioned instabilities.

Lastly, we have used the same training loop for both flows with early stopping and patience at 200, checking at every iteration, same as in [64].

Adopting the code of SPTN flows was easier, as it has been already written in Julia and the only changes we have made were in the training loop. We have used early-stopping, though due to the high cost of evaluation of validation loss, we have opted for patience 20 and checking interval 10. Note also that on higher-dimensional datasets — HAR, Letter Recognition, Isolet, Multiple Features, Arrhythmia — the implementation did not allow us to compute all 5 repetitions and therefore only the minimum of 3 has been completed.

#### D. Other reimplemented methods

As the original article on the Deep Autoencoding Gaussian Mixture [98] model (DAGMM for short), did not contain any link to a reference implementation we have been forced to follow alternative ones. We have found in total 4 similar attempts at replicating results of the original article, all of which tailored specifically to the KDDCUP99 dataset. With our implementation we bring more flexibility by implementing encoder/decoder with the same architecture as the previously described variational autoencoders on tabular data, while having much lower latent dimension  $\{1, 2\}$ . This is crucial for this method’s convergence as the latent representation together with cosine similarity and  $L_2$  norm is used to fit parameters of a Gaussian mixture, whose covariant matrices have higher chance to degenerate to non positive-semi-definitive matrices, when the latent dimension is higher than 1 on some datasets. There are two ways with which the authors and other practitioners tried to solve this problem: first is

to penalize the inverse of diagonal entries (with the weight  $\lambda_1$  of the term being on of the hyperparameters) and the second is simple addition of small  $\epsilon$  to the diagonal. Each of these remedies may in part be the cause of overall poor performance of the model, because penalizing the diagonal terms leads to mixtures with high variance and thus potentially less sensitive to the detection task. Adding one size fits all  $\epsilon$  is in this case non-systematic solution to an underlying ill posed application of this methods on particular data. We have been able to attain the performance in the original article, but the results with different models seed have high variance as has been reported in the other implementations. One other notable difference is that instead of Cholesky decomposition, we have used LU decomposition to compute covariant matrix’s determinant and it’s derivative, as it is the default in Julia’s Flux/Zygote library. We have employed early stopping based on the training loss evaluated on normal validation data for with patience of 200 iterations. Though we have seen an attempt to extend this method to image dataset, we have not pursued this direction because the problems with convergence with higher latent dimensions of the underlying autoencoder goes directly against the requirements for autoencoders on higher dimensional image data.

Situation with the REPEN [61] method, have been again easier as there exists original implementation, which we converted to Julia, while keeping the variable names and the structure itself almost intact. As with other methods, we took the liberty of extending the architecture of the input space projection to include two-layered neural networks with either relu or tanh non-linearity, which were mentioned in the article but not implemented. For its training we used fixed number of 10000 iteration of ADADelta optimizer. During prediction on test and validation data the ensemble of 1NN detector has been fitted with random subsamples of the whole training set, but with fixed seed such that the results are reproducible. This procedure has not been mentioned in the original implementation, in which training set coincided with the test set. We believe that this method would work better on contaminated training data, as the initial sorting mechanism relies on the fact that one may meaningfully split the training samples into normal data and potentially anomalous.

#### E. Evaluation

Though we have been meticulous in the implementation, due to the scale of operation we could not cover all the edge cases and therefore some runs did not complete at all or produced corrupted results with invalid values. When a method produced scores with more than half of the values being NaN, such results were filtered out. Overall less than 0.5% ( $\sim 11k$ ) of results from the tabular experiments have been affected, out of which the majority ( $\sim 10k$ ) have come from our kNN implementation on small dimensional datasets, where only low values of  $k$  produced a valid result. Results on the image datasets did not contain any invalid values.

As a result, the experimental files for some combinations of hyperparameter-seed-model were missing and had to be handled in the later stages of evaluation. In our protocol we

required at least 3 repetitions of each hyperparameter where repeated experiment were computed.

### F. Bayesian optimization

There are two parts to the implementation of Bayesian optimization:

- optimization itself and construction of parameter ranges,
- training orchestration.

The first part was mostly handled inside Python’s scikit-optimise, however the second part has proven to be quite the challenge as we have been working already with most of the result precomputed. In spite of this we believe that we have implemented quite a clean solution that is based on binary files, in which we have stored all the previously sampled/trained hyperparameters and the objective value. Each run of an experiment loaded this file, which we called Bayesian cache and fit the process from beginning with fixed seed, such that the results are reproducible. One important detail is that if an experiment did not finish, we have kept an empty entry in the cache with objective value 0.0, such that the underlying Gaussian process is given some sort of feedback. This solution allowed us to avoid long running jobs, which in some cases would take up the full 50 days or end prematurely without any mean of continuation, unless meticulously handled.

model	parameter	value set
common	dim( $\mathbf{z}$ )	{8, 16, ..., 256}
	h	{16, 32, ..., 512}
	$\eta$	{ $10^{-4}$ , $10^{-3}$ }
	batch size	{32, 64, 128}
	activation	{relu, swish, tanh}
	no. dense layers	{3, 4}
	no. conv layers	{2, 3, 4}
	channels	{16, 32, 64, 128}
	kernelsizes scalings	{(3, 5, 7, 9) (1, 2, 2, 2)}
WAE	$\lambda$	{0.1, 1}
	prior	{ $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , VampPrior}
	K	{2, 4, ..., 64}
	kernel	{rbf, img, rq}
AAE	$\sigma_k$	{ $10^{-3}$ , $10^{-2}$ , $10^{-1}$ , $10^0$ }
	$\lambda$	{0.1, 1}
	prior	{ $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , VampPrior}
	K	{2, 4, ..., 64}
	no. discriminator layers	{1, 2, 3}
adVAE	$\alpha$	{0, 0.1, ..., 1}
	dim( $\mathbf{z}$ )	{2, 4, ..., 256}
	batch size	{32, 64}
	$\gamma$	{ $5 \cdot 10^{-4}$ , $10^{-3}$ , $5 \cdot 10^{-3}$ }
	$\lambda$	{ $5 \cdot 10^{-3}$ , $10^{-2}$ , $5 \cdot 10^{-2}$ }
	$m_x$	{1, 1.5}
	$m_z$	{40, 50, 60}
(skip)GANomaly	decay	{0, 0.1, ..., 0.5}
	$w_{adv}, w_{con}, w_{enc}$	{1, 10, 20, ..., 100}
	$\lambda$	{0.1, 0.2, ..., 0.9}
	$R(x), L(X)$	{MAE, MSE}
	no. conv layers no. channels	{1, 2, 3, 4} {8, 16, ..., 128}
(fm)GAN	dim( $\mathbf{z}$ )	{2, 4, ..., 256}
	no. dense layers	{2, 3, 4}
	$\alpha$	{ $10^{-3}$ , $10^{-2}$ , ..., $10^3$ }
DeepSVDD	$\eta_{AE}$	{ $10^{-4}$ , $10^{-3}$ }
	batch norm.	{true, false}
	batch size	{64, 128}
	objective	{soft boundary, one class}
	$\nu$ decay	{0.01, 0.1, 0.5, 0.99} $10^{-6}$ }
fAnoGAN	$\eta_E$	{ $10^{-4}$ , $10^{-3}$ }
	weight clip	{0.001, 0.005, 0.01, 0.05, 0.1}
	$n_{critic}$ no. generator iterations	5 10000
DAGMM	$\eta$	{ $10^{-5}$ , $10^{-4}$ , $10^{-3}$ }
	dim( $\mathbf{z}$ )	{1, 2}
	batch size	{32, 64, ..., 256}
	no. dense layers	{2, 3}
	no. components	{2, 3, ..., 8}
	$\lambda_1$	{0.1, 0.5, 1.0}
	$\lambda_2$	{0.005, 0.05, 0.1, 0.5, 1.0}
	activation dropout	{tanh} {0.0, 0.1, ..., 0.5}
REPEN	dim( $\mathbf{z}$ )	{2, 4, ..., 64}
	h	{2, 4, ..., 256}
	confidence margin	{512, 1024, 2048}
	batch size	{32, 64, ..., 256}
	no. dense layers	{1, 2}
	ensemble size	{25, 50}
	subsample size activation	{32, 64, ..., 258} {tanh, relu}

TABLE A.4: Hyperparameters for the neural network based models.

model	parameter	value set
MAF + RealNVP	h	{16, 32, . . . , 1024}
	no. flows	{2, 4, 8}
	$\eta$	$10^{-4}$
	batch size	{32, 64, 128}
	no. layers	{2, 3}
	activation	{relu, tanh}
	batch norm.	{true, false}
	init. identity.	{true, false}
	$L_2$ reg.	{0.0, $10^{-5}$ , $10^{-6}$ }
MAF	ordering	{natural, random}
RealNVP	tanh scaling	{true, false}
SPTN	no. components	{2, 4, 8, 16}
	batch size	{32, 64, 128}
	no. flows	{1, 2, 3}
	sharing	{dense, all, none}
	first dense	{true, false}
	activation	identity

TABLE A.5: Hyperparameters for the flow-based models.

model	parameter	value
ABOD	$n$	$0.1N$
	method	fast
HBOS	no. bins	10
	$\alpha$	0.1
	tolerance	0.1
IF	no. estimators	100
	max samples	auto
	max features	1.0
kNN	$n$	$\max(10, 0.03N)$
LODA	no. bins, no. cuts	automated procedure [68]
LOF	$n$	$\max(10, 0.03N)$
OC-SVM	$\gamma$	$1/(\text{median } L_2 \text{ distance of training samples})$
	$\nu$	0.5
	kernel	rbf
PIDForest	max depth	10
	no. trees	50
	max samples	$\min(100, \max(25, N - \text{mod}(N, 50)))$
	max buckets	3
	$\epsilon$	0.1

TABLE A.6: Default hyperparameters of classical models for experiments without validation anomalies.  $N$  is the number of samples in the training dataset.

dataset	aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	osvm	pidf	maf	rnvp	sptn	fmgn	gan	mgal	dagm	dsvd	rpn	vaek	vaeo
aba	0.92	0.87	0.89	0.92	0.91	<b>0.93</b>	0.75	0.87	<b>0.93</b>	0.84	0.90	<b>0.93</b>	0.89	0.91	0.90	0.91	0.78	0.80	0.62	0.67	0.82	0.84	0.91	0.90
ann	0.82	0.83	0.80	0.84	0.88	0.78	0.89	0.78	0.78	0.69	0.80	<b>0.99</b>	0.93	0.85	0.86	0.87	0.81	0.74	0.65	0.60	0.65	0.77	0.81	0.81
arr	0.71	0.76	0.77	0.74	0.76	0.74	0.77	0.78	0.74	0.77	0.73	<b>0.81</b>	0.75	0.76	0.77	0.74	0.74	0.73	0.55	0.51	0.72	0.78	0.71	0.79
bcw	0.99	0.95	0.94	0.99	0.99	0.94	0.97	0.97	0.93	0.94	0.94	0.99	0.91	0.99	0.98	0.95	<b>1.00</b>	0.99	0.64	0.71	0.83	0.96	0.93	0.99
blt	<b>0.99</b>	0.84	0.90	0.95	0.93	0.87	0.88	0.90	0.89	0.82	0.91	0.89	0.87	0.96	0.93	0.94	0.71	0.58	0.82	0.85	0.94	0.95	0.93	0.94
bts	<b>1.00</b>	0.85	0.99	0.98	0.96	0.99	0.98	0.98	0.98	0.95	0.96	<b>1.00</b>	0.77	0.98	0.99	0.99	0.96	0.95	0.68	0.78	0.96	0.96	0.97	0.97
crd	0.88	0.61	0.69	0.72	0.53	0.56	0.50	0.69	0.61	0.74	0.67	<b>0.90</b>	0.64	0.60	0.51	0.50	0.67	0.66	0.72	0.78	0.86	0.75	0.63	0.83
eco	<b>0.90</b>	0.86	0.84	0.85	0.87	0.87	0.81	0.83	0.88	0.77	0.80	0.89	0.84	<b>0.90</b>	0.85	0.88	0.85	0.87	0.58	0.70	0.76	0.88	0.86	0.85
gls	<b>0.87</b>	0.77	0.77	0.68	0.76	0.79	0.62	0.52	0.71	0.51	0.81	0.78	0.40	0.73	0.74	0.78	0.86	0.84	0.65	0.70	0.67	0.73	0.73	0.72
hab	<b>0.97</b>	<b>0.97</b>	0.87	0.95	<b>0.97</b>	0.95	0.92	0.93	0.95	0.95	0.96	0.95	<b>0.97</b>	0.96	0.95	0.96	0.76	0.81	0.68	0.84	0.96	0.95	0.95	0.93
har	0.97	0.69	0.99	<b>1.00</b>	<b>1.00</b>	0.76	0.84	0.71	0.76	0.81	0.97	<b>1.00</b>	0.47	<b>1.00</b>	<b>1.00</b>	0.96	<b>1.00</b>	0.99	0.58	0.60	0.84	0.99	0.48	<b>1.00</b>
htr	0.96	0.93	0.94	0.96	0.96	0.95	0.96	0.95	0.95	0.95	0.95	<b>0.97</b>	0.94	0.95	0.95	0.95	0.96	0.96	0.61	0.87	0.92	0.93	0.95	0.96
ion	0.98	0.98	0.98	0.98	0.97	0.98	0.78	0.92	0.98	0.87	0.96	0.98	0.90	0.98	<b>0.99</b>	0.97	0.90	0.80	0.68	0.47	0.97	0.92	0.96	0.97
irs	0.93	0.83	0.97	0.96	0.96	0.97	0.99	0.89	0.94	<b>1.00</b>	0.88	0.93	0.99	0.79	0.80	0.93	<b>1.00</b>	0.99	0.73	0.84	0.29	0.78	0.88	0.92
iso	0.74	0.70	0.79	0.75	0.75	0.64	0.55	0.60	0.77	0.55	0.82	<b>0.84</b>	0.60	0.71	0.70	0.60	0.78	0.78	0.50	0.54	0.62	0.69	0.81	0.81
kdd	0.96	0.95	0.99	0.99	<b>1.00</b>	0.99	0.99	<b>1.00</b>	<b>1.00</b>	0.90	0.98	<b>1.00</b>	0.99	0.99	<b>1.00</b>	0.99	<b>1.00</b>	<b>1.00</b>	0.02	0.73	0.24	<b>1.00</b>	0.99	0.99
lbr	0.71	0.64	0.74	0.64	0.75	0.65	0.58	0.55	<b>0.78</b>	0.56	0.70	<b>0.78</b>	0.55	0.73	0.77	0.55	0.77	<b>0.78</b>	0.51	0.59	0.58	0.68	<b>0.78</b>	<b>0.78</b>
ltr	0.79	0.78	0.77	0.77	0.79	0.68	0.56	0.62	0.80	0.59	<b>0.83</b>	0.81	0.60	0.76	0.75	0.67	0.76	0.74	0.48	0.54	0.65	0.75	0.82	0.82
mam	0.89	0.89	0.89	0.89	0.89	0.85	0.84	0.88	0.88	0.89	0.85	<b>0.91</b>	0.86	0.87	0.89	0.88	0.78	0.82	0.75	0.89	<b>0.91</b>	0.88	0.90	0.90
mgc	0.94	0.91	0.89	<b>0.97</b>	0.96	0.94	0.83	0.90	0.94	0.82	0.93	0.94	0.91	0.96	0.96	0.96	0.85	0.84	0.55	0.52	0.81	0.85	0.89	0.90
mlt	<b>0.99</b>	0.98	0.98	0.98	0.98	0.91	0.73	0.87	0.98	0.74	0.98	<b>0.99</b>	0.83	0.98	<b>0.99</b>	0.94	<b>0.99</b>	<b>0.99</b>	0.47	0.63	0.74	0.97	<b>0.99</b>	<b>0.99</b>
mnb	0.93	0.90	0.88	0.92	0.91	0.81	0.91	0.81	0.86	0.92	0.70	<b>0.94</b>	0.82	0.90	0.90	0.86	0.73	0.83	0.67	0.65	0.85	<b>0.94</b>	0.85	0.88
pam	0.97	0.95	0.98	0.99	0.99	0.99	0.77	0.96	0.99	0.90	<b>1.00</b>	0.99	0.95	0.98	0.98	0.99	0.96	0.92	0.59	0.65	0.86	0.85	0.98	0.99
pgb	0.98	0.98	0.98	<b>0.99</b>	0.98	0.97	0.88	0.97	0.98	0.96	0.98	0.98	0.96	<b>0.99</b>	<b>0.99</b>	0.98	0.75	0.73	0.59	0.71	0.97	0.95	<b>0.99</b>	<b>0.99</b>
pim	0.83	0.78	0.81	0.86	0.86	0.83	0.81	0.83	0.84	0.81	0.82	<b>0.89</b>	0.78	0.86	0.85	0.84	0.78	0.81	0.61	0.69	0.81	0.84	0.78	0.78
prk	<b>0.89</b>	0.60	0.73	0.68	0.72	0.75	0.55	0.66	0.80	0.55	0.70	0.88	0.45	0.72	0.71	0.74	0.78	0.79	0.64	0.73	0.72	0.67	0.79	0.80
sat	0.95	0.87	0.96	0.93	0.94	0.96	0.95	0.94	0.97	0.90	0.98	<b>0.99</b>	0.95	0.91	0.93	0.84	0.97	0.97	0.74	0.77	0.82	0.97	0.95	0.97
sec	0.95	0.96	0.89	0.98	0.98	0.90	0.82	0.92	0.97	0.86	<b>0.99</b>	0.98	<b>0.99</b>	<b>0.99</b>	0.96	0.90	0.96	0.97	0.59	0.65	0.87	0.87	0.96	0.97
seg	0.92	0.91	0.94	0.90	0.92	0.95	0.86	0.90	<b>0.96</b>	0.93	0.94	0.95	0.93	0.92	0.92	0.93	0.89	0.89	0.60	0.62	0.72	0.87	0.94	0.95
sei	0.76	0.73	0.74	0.74	0.73	0.74	0.73	0.70	0.74	0.70	0.65	<b>0.77</b>	0.74	0.73	0.73	0.74	0.68	0.71	0.56	0.54	0.74	0.72	0.73	0.73
sht	0.99	0.99	0.99	<b>1.00</b>	0.98	<b>1.00</b>	0.93	0.98	<b>1.00</b>	0.90	<b>1.00</b>	<b>1.00</b>	0.99	<b>1.00</b>	0.99	<b>1.00</b>	0.85	0.87	0.65	0.76	0.93	0.99	<b>1.00</b>	<b>1.00</b>
snr	0.73	0.65	0.76	0.65	0.75	0.64	0.49	0.55	0.64	0.52	0.85	0.84	0.50	0.65	0.66	0.58	0.81	0.81	0.57	0.56	0.47	0.73	0.74	<b>0.86</b>
sph	0.80	0.35	0.52	0.69	0.28	0.35	0.30	0.35	0.50	0.47	0.40	<b>0.82</b>	0.28	0.26	0.30	0.28	0.74	0.80	0.55	0.72	0.50	0.50	0.50	0.80
spm	0.77	0.80	0.78	0.86	0.87	0.77	0.82	0.82	0.78	0.63	0.81	<b>0.94</b>	0.84	0.86	0.85	0.83	0.91	0.91	0.54	0.64	0.60	0.83	0.54	0.81
vhc	<b>0.80</b>	0.65	0.73	0.77	0.77	0.74	0.77	0.69	0.73	0.70	0.60	0.72	0.70	0.75	0.77	0.74	0.65	0.70	0.57	0.64	0.70	0.68	0.63	0.71
wf1	0.81	0.71	0.78	0.89	0.87	0.72	0.87	0.83	0.83	0.81	0.75	<b>0.95</b>	0.85	0.75	0.75	0.77	0.85	0.84	0.63	0.80	0.75	0.92	0.82	0.92
wf2	0.92	0.73	0.75	0.90	0.91	0.73	0.86	0.84	0.84	0.82	0.76	<b>0.94</b>	0.85	0.74	0.79	0.77	0.87	0.84	0.60	0.80	0.75	0.92	0.79	<b>0.94</b>
wne	<b>1.00</b>	0.98	0.96	0.99	0.98	0.95	0.92	0.91	0.98	0.83	0.97	0.99	0.73	0.98	0.95	0.96	0.97	0.92	0.61	0.66	0.94	0.93	0.95	0.99
wrb	0.73	0.73	0.81	0.86	0.85	0.80	0.87	0.81	0.82	0.72	0.76	0.85	<b>0.91</b>	0.78	0.82	0.81	0.82	0.78	0.56	0.61	0.57	0.66	0.80	0.80
yst	0.74	0.70	0.66	<b>0.75</b>	0.73	0.66	0.53	0.63	0.66	0.67	0.68	<b>0.75</b>	0.60	0.72	0.72	0.67	0.62	0.65	0.65	0.69	0.70	0.68	0.54	0.64
avg. AUC	0.88	0.81	0.85	0.87	0.86	0.82	0.78	0.81	0.85	0.78	0.84	<b>0.91</b>	0.79	0.85	0.85	0.83	0.84	0.84	0.6	0.68	0.75	0.84	0.83	0.88
avg. rank	6.9	13.4	10.3	6.5	7.0	11.4	14.9	14.5	8.4	16.1	11.4	<b>2.9</b>	14.3	8.9	9.0	10.7	11.4	12.0	22.3	19.8	16.2	12.1	11.0	6.8

TABLE B.1: Performance of models on *tabular* datasets with hyperparameter selection using 50% of all available anomalies in the validation set, reported in the AUC metric on the test data, averaged over 5 random cross-validation repetitions. The final row contains the average model rank.

dataset	aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	osvm	pidf	maf	rnvp	sptn	fmgn	gan	mgal	dagm	dsvd	rpn	vae	vaeo
aba	0.87	0.76	0.79	0.90	0.90	<b>0.93</b>	0.77	0.86	<b>0.93</b>	0.27	0.89	0.91	0.90	0.85	0.89	0.89	0.37	0.39	0.38	0.67	0.82	0.50	0.89	0.89
ann	0.82	0.58	0.69	0.79	0.76	0.67	0.66	0.74	0.67	0.43	0.64	0.64	<b>0.92</b>	0.79	0.80	0.80	0.60	0.43	0.51	0.48	0.56	0.50	0.71	0.53
arr	0.35	0.72	0.69	0.67	0.75	0.75	0.77	<b>0.78</b>	0.75	0.73	0.75	0.75	0.69	0.74	0.75	0.74	0.50	0.43	0.43	0.50	0.47	0.63	0.71	0.71
bcw	0.95	0.70	0.90	0.95	0.96	0.94	0.97	0.97	0.94	0.93	0.94	0.94	0.76	0.97	<b>0.98</b>	0.94	0.97	0.86	0.33	0.57	0.76	0.53	0.92	0.91
blt	0.67	0.88	0.87	0.89	0.88	0.91	0.90	0.94	0.93	0.68	<b>0.95</b>	0.93	0.84	0.92	0.91	0.88	0.46	0.43	0.32	0.54	0.92	0.50	0.93	0.92
bts	<b>0.99</b>	0.69	0.90	<b>0.99</b>	0.96	0.96	0.98	0.98	0.98	0.92	0.94	0.98	0.30	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	0.83	0.83	0.54	0.78	0.92	0.50	0.87	0.88
crd	0.42	0.61	0.46	0.40	0.45	0.47	0.48	0.66	0.46	0.49	0.59	0.53	0.57	0.43	0.45	0.49	<b>0.67</b>	0.51	0.50	0.45	0.62	0.50	0.52	0.56
eco	0.74	0.58	0.79	0.83	0.80	<b>0.86</b>	0.72	0.81	0.85	0.64	0.80	0.85	0.76	0.83	0.85	0.79	0.75	0.67	0.64	0.54	0.52	0.50	0.79	0.79
gls	<b>0.86</b>	0.44	0.68	0.64	0.60	0.79	0.60	0.56	0.72	0.56	0.78	0.76	0.40	0.66	0.72	0.76	0.60	0.72	0.44	0.45	0.47	0.50	0.57	0.49
hab	0.91	0.78	0.77	0.95	0.93	0.95	0.91	0.94	0.95	0.37	<b>0.96</b>	0.95	0.95	0.95	0.95	0.93	0.48	0.49	0.58	0.74	0.81	0.50	0.79	0.73
har	0.65	0.60	0.78	0.38	0.55	0.66	0.39	0.60	0.69	0.47	<b>0.93</b>	0.64	0.32	0.48	0.56	<b>0.93</b>	0.85	0.92	0.53	0.56	0.46	0.50	0.33	0.29
htr	0.93	0.86	0.94	0.94	<b>0.95</b>	-	<b>0.95</b>	<b>0.95</b>	0.94	0.35	0.94	0.93	0.92	<b>0.95</b>	0.94	<b>0.95</b>	0.79	0.93	0.26	0.84	0.91	0.50	0.94	0.94
ion	0.94	0.82	<b>0.98</b>	0.96	0.95	<b>0.98</b>	0.73	0.92	<b>0.98</b>	0.62	0.96	<b>0.98</b>	0.76	0.97	0.96	0.97	0.48	0.67	0.34	0.42	0.91	0.50	0.91	0.87
irs	0.60	0.72	0.90	0.75	0.78	0.93	0.97	0.85	0.89	0.12	0.83	0.84	<b>0.99</b>	0.75	0.80	0.93	0.95	0.97	0.59	0.68	0.04	0.50	0.88	0.64
iso	0.44	0.48	0.55	0.42	0.65	0.61	0.54	0.54	0.70	0.53	0.74	<b>0.75</b>	0.54	0.63	0.64	0.60	0.53	0.55	0.44	0.49	0.37	0.50	0.71	0.73
kdd	0.93	0.22	0.95	0.97	0.99	-	0.96	0.99	0.95	0.15	0.97	0.98	0.91	0.99	<b>1.00</b>	0.97	0.98	0.98	0.01	0.57	0.02	0.50	-	0.88
lbr	0.42	0.59	0.66	0.65	0.68	0.66	0.56	0.55	0.62	0.57	0.56	0.76	0.54	0.71	0.64	0.53	0.48	0.52	0.49	0.49	0.53	0.50	0.65	<b>0.78</b>
ltr	0.49	0.53	0.66	0.46	0.70	0.68	0.54	0.58	0.75	0.51	0.78	<b>0.79</b>	0.56	0.68	0.68	0.67	0.50	0.50	0.48	0.46	0.42	0.50	0.78	<b>0.79</b>
mam	0.74	0.78	0.87	0.84	0.81	0.85	0.85	0.87	0.87	0.18	0.83	0.86	0.83	0.84	0.87	0.87	0.36	0.48	0.31	0.67	0.88	0.50	<b>0.90</b>	0.89
mgc	0.91	0.88	0.78	0.95	0.93	0.89	0.80	0.85	0.91	0.24	0.90	0.86	0.86	<b>0.96</b>	0.95	0.95	0.74	0.80	0.49	0.41	0.70	0.50	0.84	0.82
mlt	0.49	0.55	0.76	0.45	0.96	0.89	0.73	0.80	0.96	0.63	0.96	<b>0.98</b>	0.76	0.96	0.95	0.94	0.50	0.51	0.45	0.55	0.38	0.50	0.83	0.92
mnb	0.85	0.78	0.76	0.83	0.76	-	0.81	0.80	0.85	0.12	0.81	<b>0.88</b>	0.79	0.70	0.79	0.75	0.49	0.50	0.23	0.41	0.62	0.50	0.83	0.84
pen	0.97	0.67	0.97	0.98	0.97	0.96	0.77	0.91	0.97	0.63	0.98	0.96	0.89	0.97	0.98	<b>0.99</b>	0.73	0.74	0.52	0.58	0.65	0.50	0.78	0.70
pgb	0.95	0.92	0.94	0.97	0.97	0.98	0.80	0.95	0.98	0.25	0.98	0.97	0.92	0.98	0.98	0.98	0.49	0.61	0.25	0.71	0.95	0.85	<b>0.99</b>	<b>0.99</b>
pim	0.73	0.80	0.78	0.84	0.73	<b>0.85</b>	0.82	<b>0.85</b>	0.83	0.62	0.75	0.83	0.79	0.83	0.81	0.79	0.62	0.53	0.53	0.51	0.71	0.68	0.69	0.71
prk	0.77	0.59	0.66	0.62	0.70	0.72	0.48	0.64	0.69	<b>0.81</b>	0.67	0.74	0.37	0.73	0.68	0.69	0.46	0.51	0.48	0.27	0.47	0.50	0.74	0.77
sat	0.77	0.59	0.90	0.82	0.85	0.92	0.94	0.91	0.94	0.48	0.95	0.93	0.93	0.85	0.80	0.76	0.90	<b>0.96</b>	0.74	0.55	0.66	0.50	0.92	0.89
sec	0.83	0.59	0.89	0.54	0.86	0.82	0.78	0.89	0.97	0.80	<b>0.99</b>	0.98	0.90	0.95	0.91	0.89	0.59	0.47	0.49	0.56	0.87	0.50	0.97	0.95
seg	0.83	0.39	0.87	0.88	0.88	<b>0.92</b>	0.80	0.82	<b>0.92</b>	0.28	0.89	0.91	0.87	0.86	0.86	0.88	0.80	0.79	0.44	0.43	0.19	0.50	0.91	0.90
sei	0.70	<b>0.78</b>	0.71	0.72	0.71	0.74	0.73	0.70	0.72	0.30	0.65	0.72	0.73	0.72	0.74	0.71	0.55	0.49	0.40	0.53	0.47	0.50	0.68	0.68
sht	0.99	0.82	0.97	0.98	0.97	-	0.89	0.96	0.98	0.03	0.95	0.97	0.95	<b>1.00</b>	0.99	<b>1.00</b>	0.48	0.15	0.46	0.40	0.89	0.50	0.96	0.97
snr	0.58	0.49	0.59	0.57	0.59	0.62	0.53	0.60	0.61	0.57	0.55	0.63	0.52	0.56	0.64	0.57	0.61	<b>0.65</b>	0.49	0.51	0.42	0.50	0.61	0.60
sph	<b>0.64</b>	0.34	0.42	0.49	0.30	0.25	0.25	0.29	0.23	0.55	0.22	0.24	0.21	0.24	0.23	0.25	0.53	0.51	0.46	0.48	0.32	0.50	0.25	0.28
spm	0.69	0.54	0.57	0.83	<b>0.84</b>	0.73	0.80	0.83	0.68	0.40	0.66	0.70	0.81	0.83	0.83	0.83	0.82	0.63	0.46	0.49	0.58	0.50	0.38	0.28
vhc	0.67	0.61	0.64	0.75	0.72	0.72	0.71	0.67	0.68	0.56	0.58	0.60	0.66	<b>0.76</b>	0.73	0.73	0.62	0.49	0.48	0.50	0.47	0.50	0.56	0.53
wf1	0.72	0.51	0.50	0.59	0.68	0.69	<b>0.85</b>	0.80	0.76	0.30	0.74	0.77	0.74	0.75	0.63	0.75	0.69	0.69	0.63	0.37	0.72	0.50	0.71	0.66
wf2	0.65	0.53	0.64	0.44	0.65	0.69	<b>0.85</b>	0.81	0.77	0.35	0.74	0.78	0.75	0.74	0.67	0.77	0.71	0.70	0.48	0.47	0.72	0.50	0.69	0.65
wne	0.34	0.45	0.92	0.25	0.92	0.97	0.89	0.87	0.98	0.91	0.98	<b>0.99</b>	0.56	0.91	0.88	0.93	0.93	0.85	0.61	0.66	0.94	0.50	0.94	0.96
wrb	0.73	0.48	0.59	0.81	0.75	0.70	0.72	0.69	0.67	0.44	0.66	0.70	0.74	0.78	<b>0.82</b>	0.77	0.74	0.76	0.56	0.50	0.52	0.50	0.66	0.67
yst	0.42	0.55	0.60	0.63	0.58	0.67	0.55	0.63	0.68	0.51	<b>0.70</b>	0.67	0.61	0.68	0.62	0.61	0.35	0.40	0.47	0.58	0.62	0.50	0.54	0.50
avg. AUC	0.72	0.63	0.76	0.73	0.78	0.79	0.74	0.78	<b>0.81</b>	0.49	0.80	<b>0.81</b>	0.72	0.80	0.80	0.80	0.64	0.63	0.46	0.53	0.61	0.52	0.75	0.74
avg. rank	12.5	16.2	11.6	10.4	9.0	8.6	11.3	8.6	<b>6.0</b>	18.3	8.0	<b>6.0</b>	12.5	7.2	7.1	7.4	14.8	15.4	20.6	19.7	16.6	19.5	10.7	11.6

TABLE B.2: Performance of models on *tabular* datasets with hyperparameter selection *without anomalies*. Where available, metrics akin to likelihood on validation data without anomalies were used instead. Reported values correspond to the AUC metric on the test data, averaged over 5 random cross-validation repetitions. The final row contains the average model rank. Missing values indicate a failure of a particular method to train with default parameters within the given computational budget.



(a) 50% anomalies in validation											(b) no anomalies in validation														
dataset	aae	gano	skip	vae	wae	knn	osvm	fano	fmg	dsvd	vaek	vaeo	dataset	aae	gano	skip	vae	wae	knn	osvm	fano	fmg	dsvd	vaek	vaeo
mnist:0	1.00	1.00	0.96	1.00	1.00	1.00	1.00	1.00	0.96	1.00	1.00	0.99	mnist:0	0.99	1.00	0.94	0.99	0.99	1.00	0.99	0.98	0.50	1.00	1.00	0.99
mnist:1	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00	1.00	1.00	1.00	mnist:1	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.50	1.00	1.00	1.00
mnist:2	0.96	0.97	0.86	0.96	0.96	0.93	0.93	0.94	0.89	0.98	0.94	0.94	mnist:2	0.92	0.97	0.79	0.88	0.88	0.91	0.88	0.82	0.50	0.86	0.92	0.92
mnist:3	0.97	0.97	0.87	0.97	0.97	0.96	0.94	0.95	0.74	0.97	0.97	0.96	mnist:3	0.93	0.96	0.77	0.91	0.92	0.94	0.91	0.89	0.50	0.95	0.96	0.96
mnist:4	0.98	0.97	0.92	0.98	0.97	0.96	0.96	0.96	0.85	0.98	0.97	0.96	mnist:4	0.95	0.96	0.89	0.91	0.93	0.95	0.92	0.91	0.50	0.95	0.95	0.95
mnist:5	0.98	0.98	0.85	0.98	0.98	0.97	0.96	0.96	0.78	0.97	0.97	0.97	mnist:5	0.96	0.96	0.75	0.93	0.92	0.94	0.91	0.91	0.50	0.92	0.93	0.92
mnist:6	0.99	1.00	0.96	1.00	0.99	0.99	0.98	0.99	0.93	1.00	0.99	0.99	mnist:6	0.99	0.99	0.91	0.98	0.98	0.99	0.97	0.98	0.50	0.97	0.99	0.99
mnist:7	0.98	0.99	0.95	0.98	0.98	0.97	0.96	0.97	0.96	0.98	0.97	0.96	mnist:7	0.97	0.98	0.93	0.97	0.98	0.97	0.96	0.96	0.50	0.97	0.97	0.97
mnist:8	0.92	0.96	0.81	0.92	0.92	0.91	0.90	0.94	0.80	0.97	0.91	0.92	mnist:8	0.85	0.95	0.80	0.80	0.80	0.89	0.86	0.78	0.50	0.95	0.87	0.88
mnist:9	0.98	0.98	0.93	0.98	0.98	0.97	0.96	0.98	0.91	0.99	0.98	0.98	mnist:9	0.96	0.98	0.85	0.95	0.95	0.96	0.94	0.94	0.50	0.97	0.97	0.96
boot	0.99	0.99	0.99	0.98	0.98	0.98	0.98	0.99	0.98	0.99	0.93	0.98	boot	0.98	0.99	0.97	0.98	0.98	0.98	0.97	0.97	0.50	0.95	0.89	0.87
bag	0.92	0.95	0.87	0.91	0.93	0.90	0.90	0.92	0.70	0.96	0.88	0.91	bag	0.89	0.94	0.82	0.87	0.86	0.90	0.90	0.84	0.50	0.75	0.86	0.84
coat	0.91	0.93	0.91	0.91	0.91	0.92	0.91	0.93	0.85	0.93	0.91	0.92	coat	0.90	0.93	0.90	0.89	0.87	0.91	0.91	0.89	0.50	0.86	0.90	0.89
dress	0.95	0.96	0.95	0.94	0.95	0.94	0.95	0.96	0.90	0.96	0.93	0.92	dress	0.92	0.96	0.91	0.92	0.90	0.94	0.94	0.93	0.50	0.90	0.93	0.92
pullover	0.90	0.91	0.88	0.90	0.90	0.90	0.89	0.90	0.81	0.92	0.89	0.87	pullover	0.89	0.91	0.86	0.89	0.87	0.90	0.89	0.88	0.50	0.92	0.88	0.87
sandal	0.92	0.94	0.94	0.91	0.91	0.92	0.91	0.94	0.92	0.96	0.89	0.90	sandal	0.90	0.92	0.89	0.89	0.90	0.90	0.90	0.91	0.50	0.74	0.83	0.84
shirt	0.84	0.87	0.81	0.84	0.84	0.85	0.84	0.86	0.59	0.87	0.85	0.83	shirt	0.82	0.87	0.78	0.81	0.78	0.85	0.84	0.80	0.50	0.78	0.84	0.82
sneaker	0.98	0.99	0.99	0.99	0.99	0.98	0.98	0.99	0.96	0.99	0.96	0.98	sneaker	0.98	0.99	0.98	0.98	0.98	0.98	0.98	0.98	0.50	0.73	0.95	0.95
t-shirt	0.93	0.94	0.91	0.93	0.93	0.93	0.93	0.94	0.72	0.95	0.92	0.93	t-shirt	0.91	0.94	0.89	0.90	0.88	0.93	0.93	0.87	0.50	0.92	0.90	0.89
trouser	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.97	0.99	0.99	0.99	trouser	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.50	0.99	0.98	0.99
bright	1.00	0.93	0.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	bright	1.00	0.00	0.00	1.00	1.00	1.00	1.00	1.00	0.50	0.22	0.95	0.98
cannye	1.00	0.84	0.63	1.00	1.00	1.00	1.00	1.00	0.98	0.99	1.00	1.00	cannye	1.00	0.28	0.12	1.00	1.00	0.99	0.97	1.00	0.50	0.08	1.00	1.00
dottedl	0.99	0.64	0.43	1.00	1.00	0.92	0.88	1.00	0.98	0.89	0.96	0.92	dottedl	0.99	0.51	0.28	1.00	1.00	0.80	0.82	1.00	0.50	0.59	0.88	0.86
fog	1.00	0.15	0.00	1.00	1.00	0.99	0.97	1.00	1.00	1.00	0.99	1.00	fog	1.00	0.00	0.00	1.00	1.00	0.92	0.80	1.00	0.56	0.08	0.95	0.93
glassb	0.97	0.45	0.14	0.98	1.00	0.76	0.71	0.97	1.00	0.58	0.71	1.00	glassb	0.97	0.03	0.00	0.97	0.97	0.26	0.12	0.96	0.51	0.02	0.30	0.21
impulsn	1.00	0.78	0.11	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	impulsn	1.00	0.58	0.00	1.00	1.00	1.00	1.00	1.00	0.52	0.12	1.00	1.00
motionb	0.80	0.44	0.20	0.80	0.98	0.72	0.71	0.99	1.00	0.63	0.90	0.99	motionb	0.37	0.19	0.00	0.71	0.95	0.29	0.14	0.50	0.54	0.36	0.41	0.33
rotate	0.61	0.51	0.37	0.65	0.64	0.65	0.64	0.68	0.95	0.59	0.68	0.87	rotate	0.24	0.35	0.12	0.51	0.63	0.50	0.45	0.43	0.50	0.44	0.49	0.48
scale	0.53	0.88	0.91	0.58	0.59	0.57	0.95	0.76	0.94	0.79	0.91	0.99	scale	0.21	0.12	0.65	0.22	0.15	0.30	0.11	0.17	0.50	0.09	0.40	0.33
shear	0.64	0.55	0.28	0.83	0.93	0.65	0.64	0.77	0.97	0.72	0.62	0.84	shear	0.57	0.40	0.08	0.83	0.93	0.58	0.54	0.54	0.50	0.47	0.55	0.54
shotn	0.99	0.58	0.56	1.00	1.00	0.84	0.73	0.99	0.87	0.58	0.91	0.96	shotn	0.99	0.34	0.21	1.00	1.00	0.64	0.62	0.94	0.59	0.07	0.78	0.77
spatter	0.96	0.46	0.22	0.97	1.00	0.70	0.55	0.91	1.00	0.68	0.73	0.93	spatter	0.96	0.23	0.01	0.96	1.00	0.57	0.56	0.91	0.97	0.40	0.62	0.59
stripe	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	stripe	1.00	0.33	0.95	1.00	1.00	1.00	1.00	1.00	0.53	0.74	1.00	1.00
transit	0.99	0.70	0.62	1.00	0.99	0.99	0.97	0.97	0.72	0.90	0.99	0.99	transit	0.98	0.55	0.51	0.95	0.96	0.94	0.92	0.84	0.52	0.70	0.97	0.96
grid	0.50	0.62	0.62	0.60	0.61	0.48	0.53	0.70	0.67	0.60	0.38	0.61	grid	0.53	0.50	0.57	0.60	0.50	0.45	0.49	0.40	0.50	0.51	0.37	0.40
transtr	0.79	0.74	0.69	0.81	0.80	0.76	0.75	0.78	0.77	0.70	0.74	0.73	transtr	0.78	0.56	0.68	0.81	0.78	0.75	0.72	0.64	0.50	0.69	0.74	0.75
wood	0.77	0.57	0.54	0.77	0.75	0.77	0.75	0.69	0.71	0.56	0.69	0.65	wood	0.77	0.57	0.52	0.77	0.76	0.75	0.74	0.59	0.50	0.54	0.64	0.66
avg. AUC	0.91	0.81	0.69	0.92	0.93	0.89	0.88	0.93	0.89	0.88	0.89	0.93	avg. AUC	0.87	0.67	0.6	0.89	0.89	0.83	0.8	0.84	0.52	0.65	0.83	0.82
avg. rank	3.6	5.1	9.1	3.1	2.9	5.4	6.4	3.4	7.5	4.0	5.7	4.9	avg. rank	3.2	5.2	9.1	3.6	4.1	3.8	5.7	6.0	10.1	7.5	4.7	5.2

(a) 50% anomalies in validation

(b) no anomalies in validation

TABLE B.3: Performance of models on the image data containing statistical anomalies with hyperparameter selection based on the availability of anomalies in validation data. Reported are the AUC values on the test dataset, averaged 5 random cross-validation repetitions where available.

(a) 50% anomalies in validation											(b) no anomalies in validation														
dataset	aae	gano	skip	vae	wae	knn	osvm	fano	fmgn	dsvd	vae	vaeo	dataset	aae	gano	skip	vae	wae	knn	osvm	fano	fmgn	dsvd	vae	vaeo
airpln	0.69	0.65	0.76	0.68	0.68	0.70	0.72	0.77	<b>0.81</b>	0.75	0.68	0.68	airpln	0.67	0.52	0.53	0.68	0.68	0.66	0.68	0.67	0.50	<b>0.69</b>	0.64	0.65
automb	0.43	0.54	0.67	0.62	0.52	0.46	0.52	0.60	<b>0.83</b>	0.67	0.51	0.60	automb	0.35	0.50	0.46	0.36	0.35	0.42	0.45	0.37	0.50	<b>0.67</b>	0.42	0.43
bird	0.69	0.68	0.68	0.68	0.69	0.70	0.70	0.69	0.69	0.67	<b>0.71</b>	0.70	bird	0.69	0.62	0.65	0.68	0.68	0.69	0.69	0.68	0.50	0.58	0.69	<b>0.70</b>
cat	0.56	0.54	<b>0.66</b>	0.59	0.60	0.52	0.51	0.58	0.61	0.61	0.51	0.51	cat	0.56	0.50	<b>0.60</b>	0.59	<b>0.60</b>	0.50	0.50	0.58	0.50	0.50	0.48	0.50
deer	0.76	<b>0.77</b>	0.75	0.75	0.73	0.76	0.76	0.75	0.60	0.74	0.76	0.76	deer	0.72	0.73	0.70	0.70	0.69	0.75	<b>0.76</b>	0.74	0.50	0.72	0.74	0.75
dog	0.57	0.56	<b>0.67</b>	0.60	0.61	0.53	0.53	0.60	0.66	0.60	0.52	0.54	dog	0.57	0.48	0.57	0.60	<b>0.61</b>	0.52	0.52	0.54	0.50	0.53	0.49	0.50
frog	0.74	<b>0.76</b>	0.74	<b>0.76</b>	0.70	0.73	0.73	<b>0.76</b>	0.70	0.72	<b>0.76</b>	0.75	frog	0.62	<b>0.73</b>	0.72	0.56	0.52	0.72	<b>0.73</b>	0.64	0.50	0.50	0.72	0.72
horse	0.54	0.58	0.64	0.58	0.54	0.54	0.54	0.60	<b>0.76</b>	0.61	0.56	0.56	horse	0.50	0.50	0.50	0.50	0.51	0.52	<b>0.53</b>	0.49	0.50	<b>0.53</b>	0.49	0.49
ship	0.72	0.67	0.78	0.72	0.72	0.71	0.70	<b>0.80</b>	0.76	0.74	0.70	0.69	ship	0.69	0.57	0.61	0.71	<b>0.72</b>	0.69	0.70	<b>0.72</b>	0.50	0.68	0.66	0.67
truck	0.54	0.54	0.66	0.66	0.55	0.48	0.59	0.66	<b>0.76</b>	0.70	0.46	0.62	truck	0.36	0.45	0.40	0.36	0.35	0.43	0.46	0.36	0.50	<b>0.64</b>	0.41	0.42
svhn2:0	0.65	0.61	0.54	0.65	0.65	0.61	0.60	0.66	<b>0.80</b>	0.65	0.64	0.62	svhn2:0	0.59	0.61	0.50	0.56	0.56	0.58	0.57	0.55	0.50	<b>0.63</b>	0.60	0.59
svhn2:1	0.68	0.61	0.57	0.68	0.68	0.63	0.61	0.62	<b>0.79</b>	0.63	0.67	0.64	svhn2:1	0.62	0.60	0.54	0.59	0.58	0.58	0.57	0.62	0.50	<b>0.63</b>	0.62	0.61
svhn2:2	0.62	0.58	0.53	0.62	0.64	0.59	0.58	0.58	<b>0.74</b>	0.61	0.62	0.59	svhn2:2	0.56	0.57	0.51	0.53	0.53	0.55	0.55	0.53	0.50	<b>0.58</b>	<b>0.58</b>	0.57
svhn2:3	0.59	0.56	0.52	0.59	0.58	0.57	0.56	0.57	<b>0.68</b>	0.58	0.59	0.58	svhn2:3	0.52	<b>0.55</b>	0.49	0.50	0.50	0.54	0.54	0.52	0.50	0.54	<b>0.55</b>	0.54
svhn2:4	0.63	0.59	0.52	0.65	0.66	0.61	0.59	0.62	<b>0.74</b>	0.58	0.64	0.61	svhn2:4	0.57	0.57	0.52	0.54	0.54	0.57	0.56	0.54	0.50	0.58	<b>0.59</b>	0.58
svhn2:5	0.60	0.57	0.52	0.60	0.61	0.57	0.57	0.59	<b>0.69</b>	0.57	0.60	0.59	svhn2:5	0.53	0.49	0.49	0.51	0.50	0.55	0.55	0.52	0.50	0.55	<b>0.56</b>	<b>0.56</b>
svhn2:6	0.59	0.56	0.53	0.59	0.59	0.58	0.56	0.60	<b>0.72</b>	0.57	0.60	0.56	svhn2:6	0.52	<b>0.56</b>	0.49	0.50	0.50	0.55	0.54	0.49	0.50	<b>0.56</b>	0.52	0.51
svhn2:7	0.65	0.60	0.53	0.66	0.66	0.61	0.60	0.61	<b>0.80</b>	0.61	0.64	0.62	svhn2:7	<b>0.61</b>	0.59	0.51	0.57	0.56	0.57	0.56	0.56	0.50	0.60	0.60	0.59
svhn2:8	0.60	0.58	0.52	0.60	0.60	0.56	0.55	0.60	<b>0.74</b>	0.57	0.59	0.58	svhn2:8	0.52	0.54	0.47	0.49	0.49	0.53	0.53	0.53	0.50	<b>0.56</b>	0.55	0.54
svhn2:9	0.60	0.57	0.51	0.60	0.60	0.57	0.57	0.59	<b>0.70</b>	0.58	0.60	0.58	svhn2:9	0.53	0.56	0.48	0.51	0.51	0.54	0.53	0.54	0.50	<b>0.57</b>	0.56	0.55
avg. AUC	0.62	0.61	0.62	0.64	0.63	0.6	0.6	0.64	<b>0.73</b>	0.64	0.62	0.62	avg. AUC	0.56	0.56	0.54	0.55	0.55	0.57	0.58	0.56	0.5	<b>0.59</b>	0.57	0.57
avg. rank	5.0	8.1	7.8	3.8	4.9	7.8	8.3	4.8	<b>2.4</b>	6.2	5.6	6.6	avg. rank	5.4	5.2	8.8	7.0	7.2	4.8	4.6	6.6	9.4	<b>3.4</b>	4.7	4.8

(a) 50% anomalies in validation

(b) no anomalies in validation

TABLE B.4: Performance of models on the image data containing semantic anomalies with hyperparameter selection based on the availability of anomalies in validation data. Reported are the AUC values on the test dataset, averaged 5 random cross-validation repetitions where available.

aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	orbf	pidf	maf	rnpv	sptn	fmgn	gan	mgal	dagm	dsvd	rpn	vae	vaeo	
avg. rank	6.6	13.1	10.0	<b>6.2</b>	6.7	11.2	14.6	14.3	8.0	15.8	11.2	8.1	14.2	8.6	8.6	10.4	11.1	11.7	22.3	19.6	16.1	11.8	10.8	6.4

TABLE D.1: Average AUC model ranks on tabular data, where **osvm** from Tab. B.1 has been replaced with **orbf** - optimizing the width of rbf kernel and using  $\nu = 0.5$ .

metric	aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	osvm	pidf	maf	rnpv	sptn	fmgn	gan	mgal	dagm	dsvd	rpn	vae	vaeo
AUC	6.9	13.4	10.3	6.5	7.0	11.4	14.9	14.5	8.4	16.1	11.4	<b>2.9</b>	14.3	8.9	9.0	10.7	11.4	12.0	22.3	19.8	16.2	12.1	11.0	6.8
TPR@5	7.0	12.6	10.1	8.4	9.6	12.3	15.1	16.2	10.2	18.0	11.1	<b>3.2</b>	14.1	10.1	10.2	12.2	9.7	10.1	19.2	20.6	16.1	12.0	11.5	6.6
rank. change	0.1	-0.8	-0.2	1.9	<b>2.6</b>	0.9	0.2	1.7	1.8	1.9	-0.3	0.3	-0.2	1.2	1.2	1.5	-1.7	-1.9	<b>-3.1</b>	0.8	-0.1	-0.1	0.5	-0.2

(a) tabular

metric	aae	gano	skip	vae	wae	knn	osvm	fano	fmgn	dsvd	vae	vaeo
AUC	3.6	5.1	9.1	3.1	<b>2.9</b>	5.4	6.4	3.4	7.5	4.0	5.7	4.9
TPR@5	3.9	5.2	9.3	<b>3.8</b>	<b>3.8</b>	6.7	6.7	4.7	7.5	4.2	6.7	6.1
rank. change	0.3	0.1	0.2	0.7	0.9	1.3	0.3	<b>1.3</b>	<b>0.0</b>	0.2	1.0	1.2

(b) statistic

metric	aae	gano	skip	vae	wae	knn	osvm	fano	fmgn	dsvd	vae	vaeo
AUC	5.0	8.1	7.8	3.8	4.9	7.8	8.3	4.8	<b>2.4</b>	6.2	5.6	6.6
TPR@5	5.8	8.2	7.0	4.8	4.1	7.8	8.8	4.2	<b>2.6</b>	5.4	6.6	6.3
rank. change	0.8	0.1	-0.8	<b>1.0</b>	<b>-0.8</b>	0.0	0.5	-0.6	0.2	-0.8	<b>1.0</b>	-0.3

(c) semantic

TABLE E.1: Comparison of average ranks on different dataset types in the anomaly validation context for two selection criteria: AUC or TPR@5. The results of hyperparameters chosen according to each metric have been averaged over 5 random cross-validation repetitions, where available and reported ranks are averaged over all datasets in the respective categories. The last row shows the rank difference between both metrics.

Model	Examples										AUC
blpix											0.68
vae											0.68
vaeo											0.68
vaek											0.68
aae											0.69
wae											0.68
gano											0.65
skip											0.76
osvm											0.72
knn											0.70
fmgn											0.81
fano											0.77
dsvd											0.75

Fig. C.1: CIFAR10 images with the lowest anomaly scores (*most normal*) for each model. Models was trained only on airplane images. Blpix is detector based on number of blue pixels.

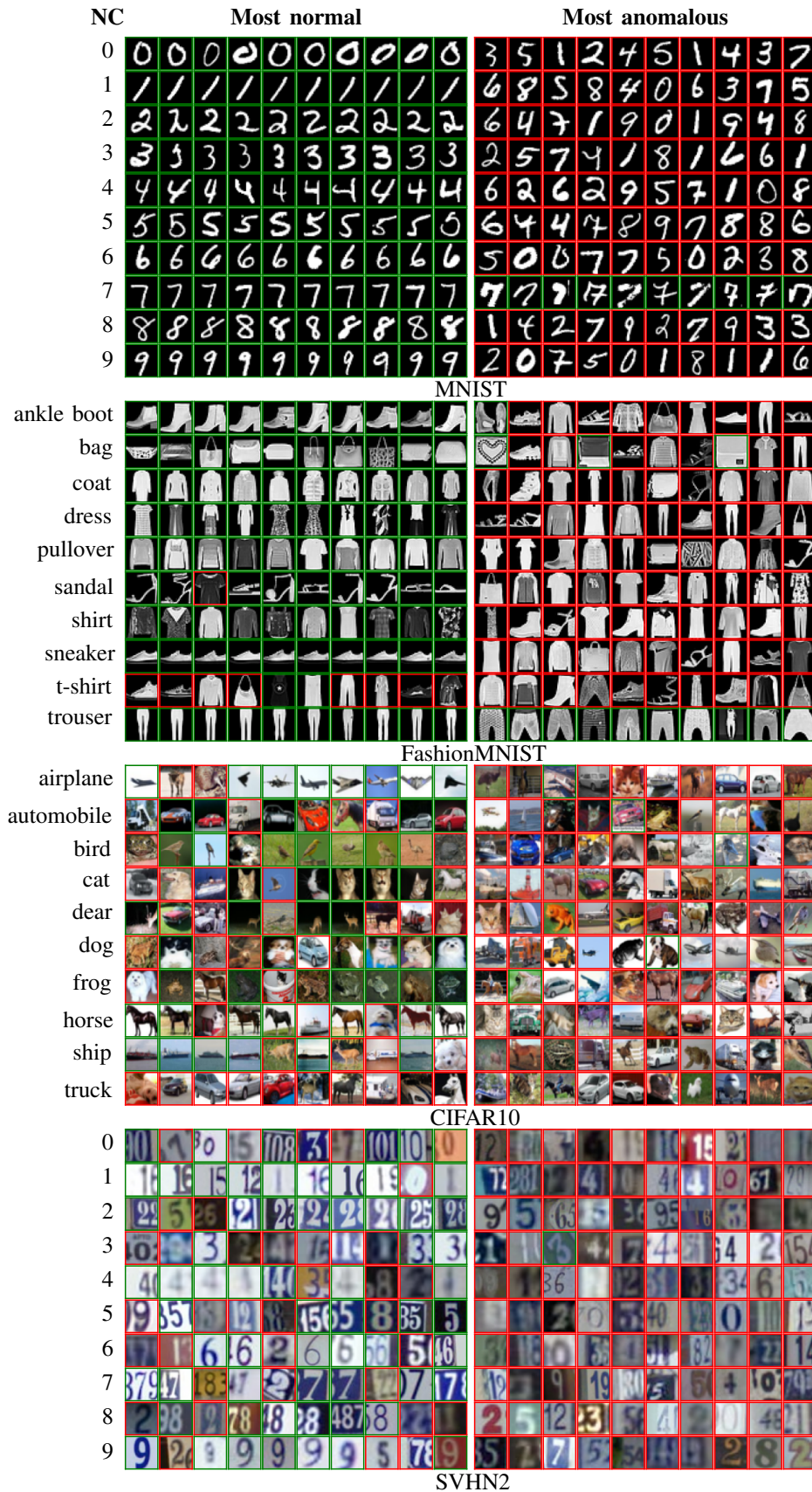


Fig. C.2: Examples of fmGAN (fmgn), NC means normal class on which the models were trained. Color of frame around images corresponds to their true class (labels): *green* means normal and *red* means anomalous.

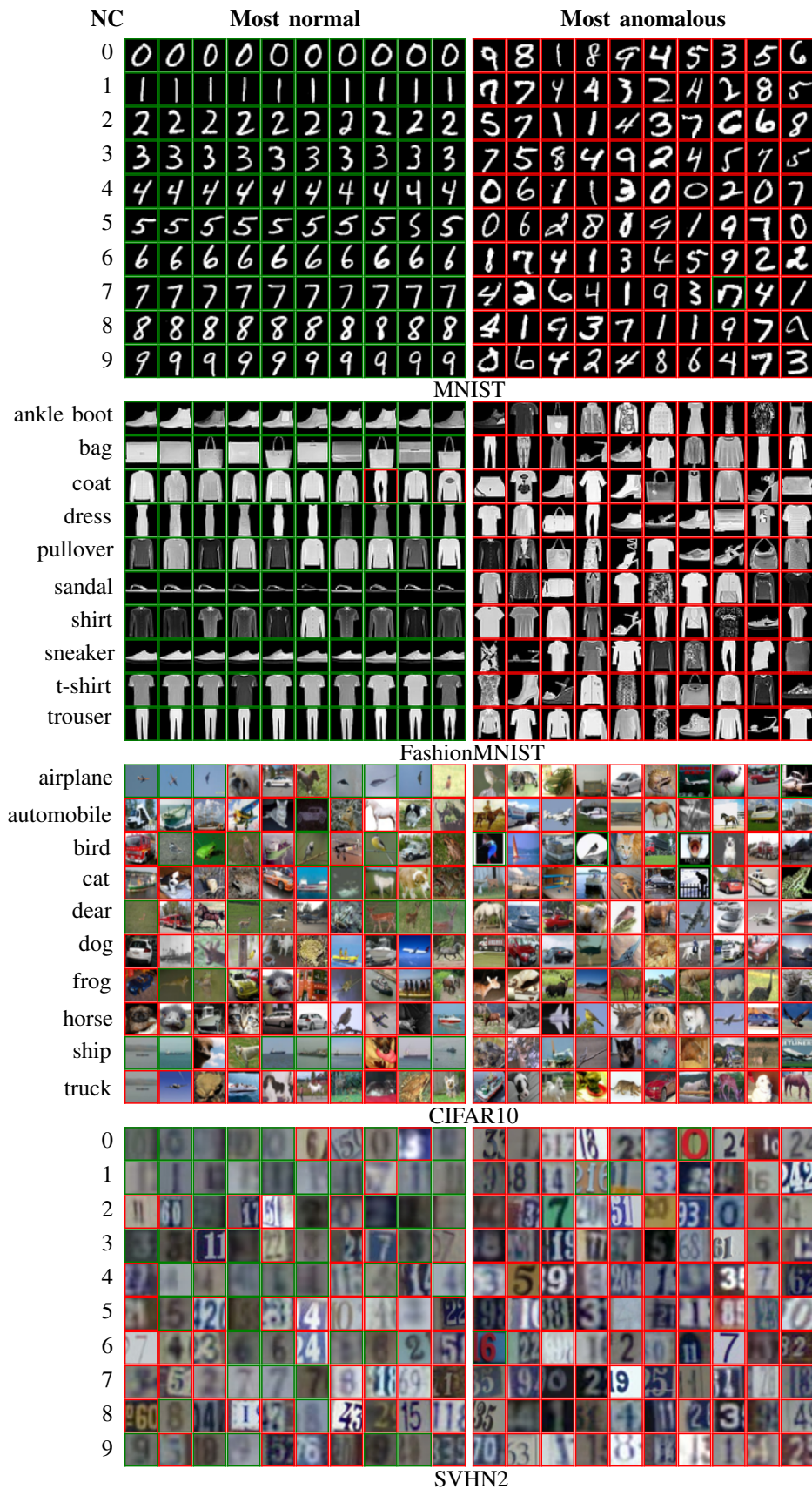


Fig. C.3: Examples of WAE with, NC means normal class on which the models were trained. Color of frame around images corresponds to their true class (labels): *green* means normal and *red* means anomalous.

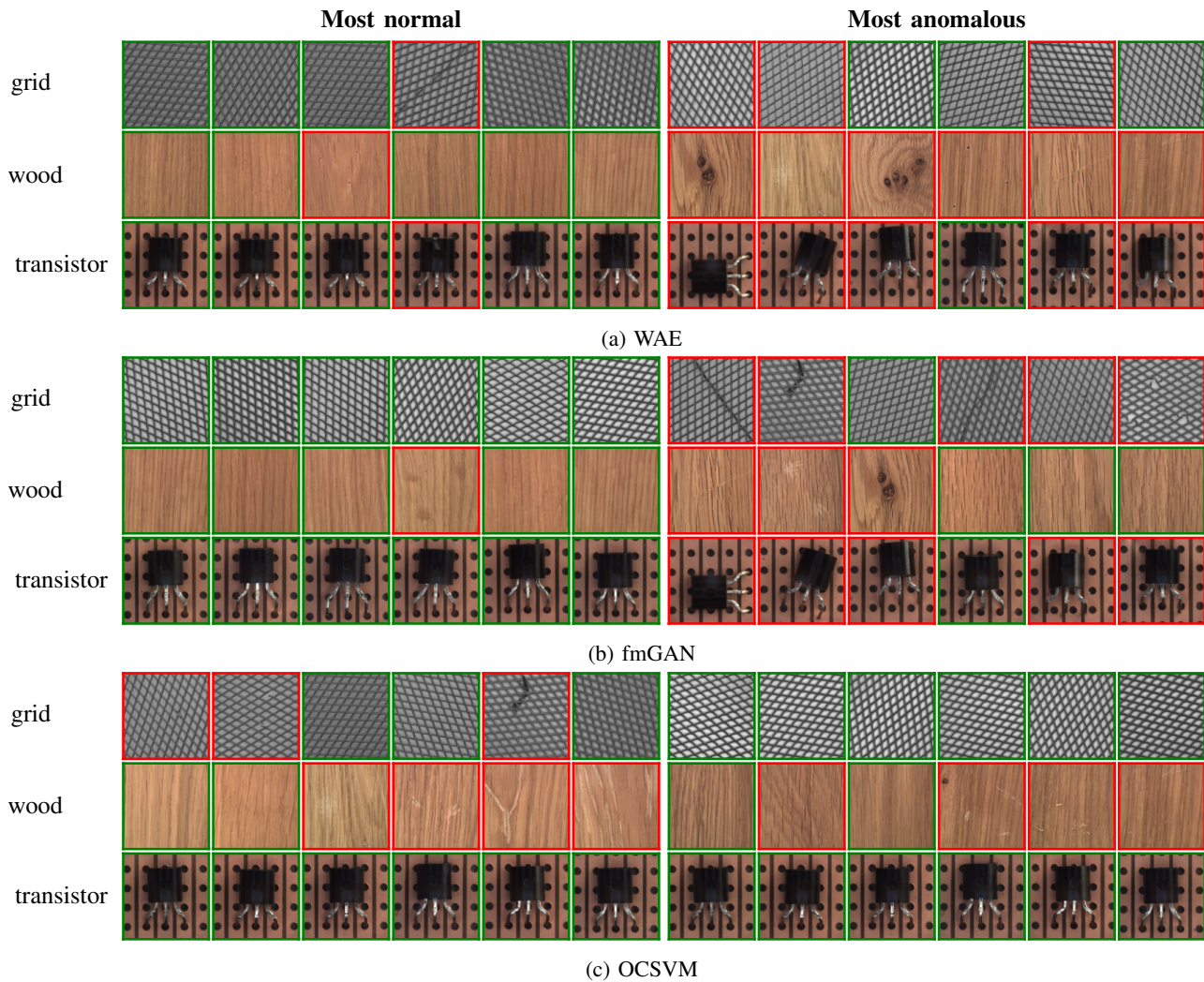
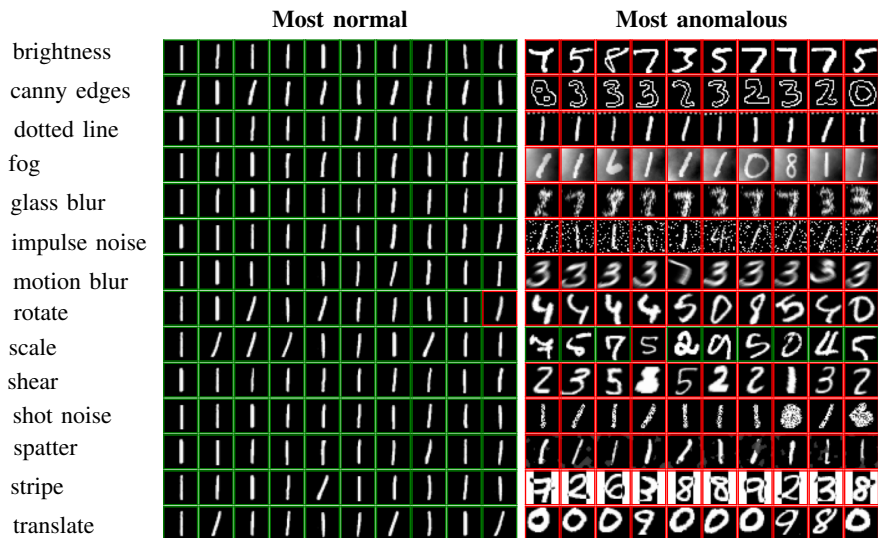
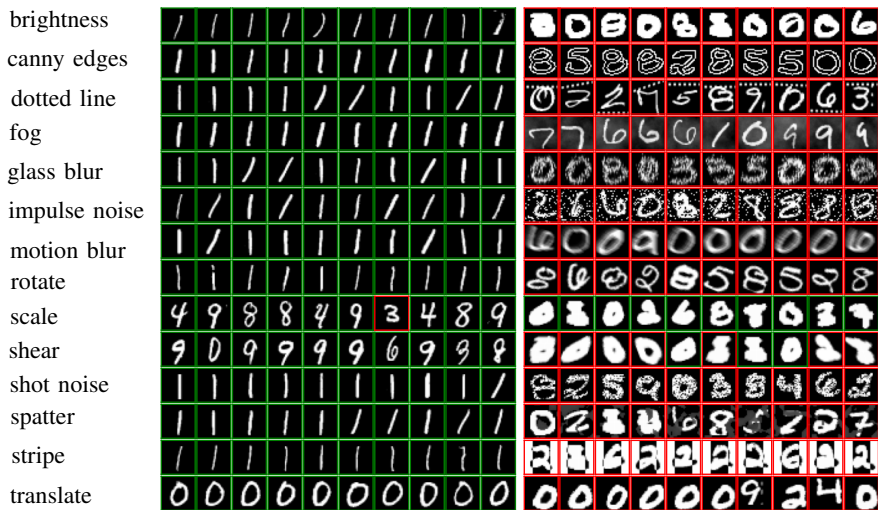


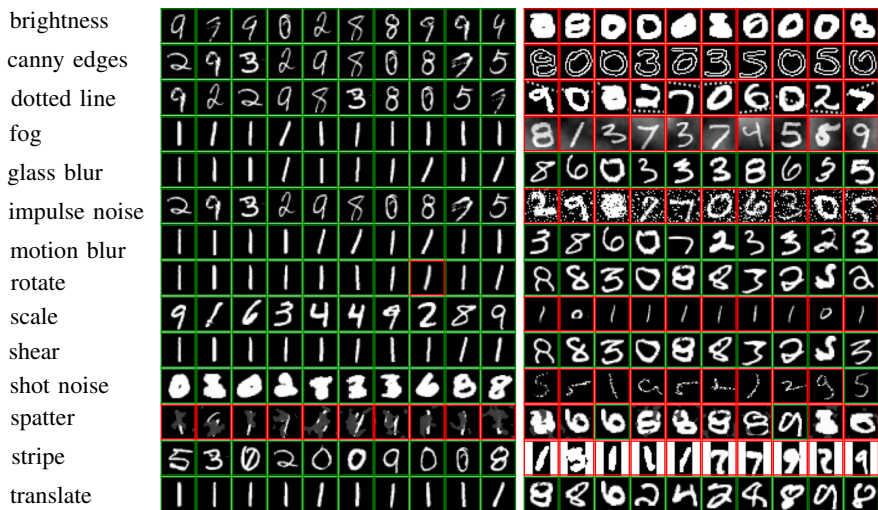
Fig. C.4: Comparison of models on MVTEC-AD dataset. Left column consists of samples with the lowest anomaly score (most normal) and right column with the highest score (most anomalous). Color of frame around images corresponds to their true class (labels): *green* means normal and *red* means anomalous.



(a) WAE



(b) fAnoGAN



(c) OCSVM

Fig. C.5: Comparison of models on MNIST-C dataset. Left column consists of samples with the lowest anomaly score (most normal), right column with the highest score (most anomalous) and each row is different anomalous corruption. Color of frame around images corresponds to their true class (labels): *green* means normal and *red* means anomalous.

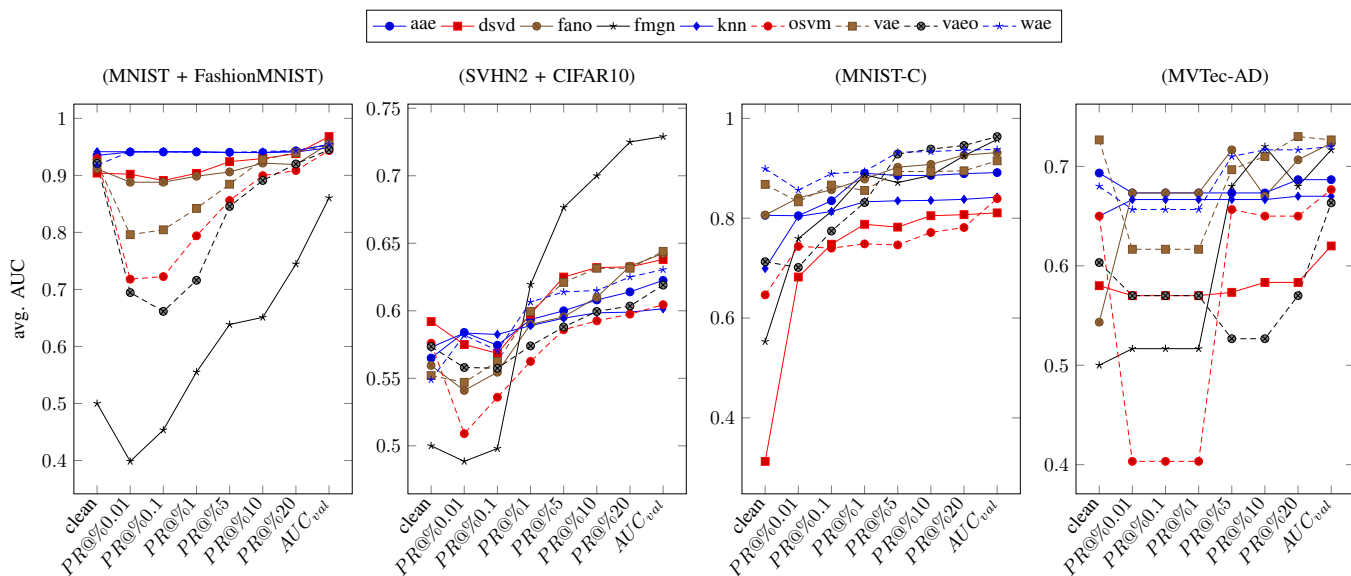


Fig. D.1: Sensitivity of methods to the number of anomalies available in the validation set for hyperparameter selection visualized in terms of the achieved AUC aggregated over all datasets in each category (columns). The clean validation context is the left-most point on the x-axis, and the anomaly validation context (50% of available anomalies) is the right-most point. The points in-between were obtained by selecting models with highest precision on the reported portion (e.g. 5%) of validation samples with the highest anomaly scores.



<b>parameter</b>	value
estimator	GP with Matern kernel
acq. function	automatic choice
acq. function optimizer	auto
no. restarts of optimizer	5
acq. function no. points	10000
noise	gaussian
$\xi$	0.01
$\kappa$	1.96

TABLE E.2: Default parameters of bayesian optimization from scikit-optimise [65] framework used for all models and datasets.

	aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	osvm	pidf	maf	rnvp	sptn	fmg	gan	mgal	dagm	dsvd	rpn	vaek	vaeo
random sampl.	6.9	13.4	10.3	6.5	7.0	11.4	14.9	14.5	8.4	16.1	11.4	<b>2.9</b>	14.3	8.9	9.0	10.7	11.4	12.0	22.3	19.8	16.2	12.1	11.0	6.8
Bayesian opt.	7.7	15.6	7.0	5.6	6.2	12.3	14.9	15.3	9.4	16.5	12.3	<b>2.8</b>	14.8	7.9	6.9	11.1	11.0	9.6	22.4	19.5	16.2	11.6	11.8	7.3
rank change	0.8	2.2	<b>-3.3</b>	-0.9	-0.8	0.9	0.0	0.8	1.0	0.4	0.9	-0.1	0.5	-1.0	-2.1	0.4	-0.4	-2.4	0.1	-0.3	0.0	-0.5	0.8	0.5
avg. change	-0.00	-0.02	0.03	0.01	0.00	0.00	0.01	-0.00	0.00	0.01	0.00	0.01	0.00	0.01	0.02	0.00	0.01	0.02	0.00	<b>0.04</b>	0.03	0.01	0.00	-0.00
no. improved	9	14	<b>30</b>	15	14	0	15	3	0	15	0	8	7	16	25	16	23	27	0	18	26	17	4	3
no. degraded	13	<b>21</b>	3	5	6	0	7	10	0	7	0	4	5	4	2	3	7	2	0	9	6	5	0	2

TABLE E.3: Comparison of average ranks in the AUC metric between Bayesian sampling and random search sampling. Methods that were not optimized with Bayesian are marked with magenta color. The second to last and last row shows the number of datasets on which the performance improved and degraded respectively.

crit	size	aae	avae	gano	vae	wae	abod	hbos	if	knn	loda	lof	osvm	pidf	maf	rnvp	sptn	fmg	gan	mgal	dagm	dsvd	rpn	vaek	vaeo
AUC	5	7.8	17.8	7.2	6.5	6.4	13.0	16.4	14.5	9.8	15.6	11.8	<b>4.0</b>	14.0	7.6	7.8	11.3	10.2	9.7	17.2	19.8	17.8	12.6	11.0	6.4
		0.01	-0.04	<b>0.02</b>	0.00	0.00	-0.00	-0.01	-0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.01	0.01	-0.00	-0.02	-0.01	0.00	-0.00
	10	7.3	19.8	7.5	6.3	5.9	12.6	16.1	14.2	10.2	15.3	11.8	<b>4.1</b>	13.9	7.6	7.9	11.2	9.7	9.4	17.9	19.8	18.5	12.5	10.8	6.9
		0.00	-0.11	<b>0.01</b>	0.00	0.00	-0.00	-0.02	-0.01	-0.01	-0.01	-0.00	-0.00	-0.00	0.00	0.00	0.00	0.00	0.00	-0.04	-0.02	-0.04	-0.01	0.00	-0.01
TPR@5	5	8.1	18.3	7.0	6.7	6.6	11.4	15.4	12.9	9.1	15.4	11.0	<b>4.8</b>	14.0	8.0	7.6	11.4	10.0	10.4	17.2	22.2	18.4	12.8	10.7	8.7
		-0.01	-0.07	<b>-0.00</b>	-0.01	-0.01	<b>-0.00</b>	-0.03	-0.01	-0.01	-0.02	-0.01	-0.01	-0.02	-0.01	-0.01	-0.01	-0.01	-0.01	-0.04	-0.17	-0.06	-0.03	-0.02	-0.04
	10	7.8	19.2	7.5	6.7	6.3	11.0	15.3	13.1	9.2	15.4	11.0	<b>4.7</b>	13.9	7.8	8.0	11.0	9.8	10.0	18.4	22.4	18.6	12.8	10.3	8.9
		-0.01	-0.14	<b>-0.00</b>	-0.01	-0.01	<b>-0.00</b>	-0.03	-0.01	-0.02	-0.03	-0.01	-0.02	-0.02	-0.01	-0.01	-0.01	-0.02	-0.02	-0.12	-0.24	-0.08	-0.03	-0.02	-0.04

(a) tabular

crit	size	aae	gano	skip	vae	wae	knn	osvm	fano	fmg	dsvd	vaek	vaeo
AUC	5	3.5	5.0	9.1	3.2	<b>3.1</b>	5.4	7.0	3.4	7.1	4.5	5.7	6.5
		<b>-0.00</b>	<b>-0.00</b>	-0.01	<b>-0.00</b>	<b>-0.00</b>	-0.01	-0.01	<b>-0.00</b>	-0.03	<b>-0.00</b>	-0.01	
	10	3.1	4.9	8.9	3.0	<b>2.9</b>	5.2	8.1	3.9	6.9	4.9	5.6	6.2
		-0.01	-0.02	-0.02	<b>-0.00</b>	<b>-0.00</b>	-0.03	-0.02	-0.01	-0.05	<b>-0.00</b>	-0.02	
TPR@5	5	3.7	4.9	9.3	3.1	<b>2.9</b>	5.2	7.2	3.8	7.2	4.8	5.4	6.1
		-0.01	-0.01	-0.03	-0.01	<b>-0.00</b>	<b>-0.00</b>	-0.03	-0.01	-0.01	-0.03	-0.01	-0.01
	10	3.4	4.9	9.2	3.1	<b>2.9</b>	5.2	7.9	4.2	7.1	4.8	5.4	6.4
		-0.02	-0.03	-0.03	<b>-0.01</b>	<b>-0.01</b>	<b>-0.01</b>	-0.04	-0.02	-0.02	-0.05	<b>-0.01</b>	-0.02

(b) statistic

crit	size	aae	gano	skip	vae	wae	knn	osvm	fano	fmg	dsvd	vaek	vaeo
AUC	5	5.1	7.6	7.6	3.4	5.0	7.6	8.3	5.2	<b>2.8</b>	6.3	5.3	7.4
		-0.00	-0.00	0.00	-0.00	-0.01	-0.00	-0.00	-0.01	<b>0.02</b>	0.00	-0.00	-0.01
	10	5.0	7.4	7.2	3.4	5.2	7.2	8.8	5.7	<b>2.8</b>	5.6	5.4	8.4
		-0.00	-0.01	0.00	-0.00	-0.01	-0.00	-0.01	-0.02	<b>0.01</b>	-0.00	-0.01	-0.02
TPR@5	5	5.6	7.8	8.1	4.5	4.8	6.7	8.6	5.6	<b>2.6</b>	5.8	5.0	8.0
		-0.01	-0.01	-0.01	-0.02	-0.01	-0.00	-0.02	-0.03	<b>0.01</b>	-0.00	-0.01	-0.03
	10	5.0	7.4	7.7	4.4	4.8	6.8	9.8	5.2	<b>2.7</b>	5.2	5.0	8.4
		-0.01	-0.02	-0.01	-0.03	-0.02	-0.01	-0.03	-0.03	<b>0.01</b>	-0.00	-0.01	-0.04

(c) semantic

TABLE E.4: Comparison of the average AUC ranks and average change in the AUC with different ensemble sizes and criterion used for their selection for different types of data. Rows with values in blue show average improvement in AUC from baseline - no ensemble setting.