

ARES: Locally Adaptive Reconstruction-based Anomaly Scoring

Adam Goodge^{1,3}, Bryan Hooi^{1,2}, See Kiong Ng^{1,2}, and Wee Siong Ng³

¹ School of Computing, National University of Singapore

² Institute of Data Science, National University of Singapore

³ Institute for Infocomm Research, A*STAR, Singapore

adam.goodge@u.nus.edu, {dcsbkh, seekiong}@nus.edu.sg,
wsng@i2r.a-star.edu.sg

Abstract. How can we detect anomalies: that is, samples that significantly differ from a given set of high-dimensional data, such as images or sensor data? This is a practical problem with numerous applications and is also relevant to the goal of making learning algorithms more robust to unexpected inputs. Autoencoders are a popular approach, partly due to their simplicity and their ability to perform dimension reduction. However, the anomaly scoring function is not adaptive to the natural variation in reconstruction error across the range of normal samples, which hinders their ability to detect real anomalies. In this paper, we empirically demonstrate the importance of local adaptivity for anomaly scoring in experiments with real data. We then propose our novel Adaptive Reconstruction Error-based Scoring approach, which adapts its scoring based on the local behaviour of reconstruction error over the latent space. We show that this improves anomaly detection performance over relevant baselines in a wide variety of benchmark datasets.

Keywords: Anomaly detection · machine learning · unsupervised learning.

1 Introduction

The detection of anomalous data is important in a wide variety of applications, such as detecting fraudulent financial transactions or malignant tumours. Recently, deep learning methods have enabled significant improvements in performance on ever larger and higher-dimensional datasets. Despite this, anomaly detection remains a challenging task, most notably due to the difficulty of obtaining accurate labels of anomalous data. For this reason, supervised classification methods are unsuitable for anomaly detection. Instead, unsupervised methods are used to learn the distribution of an all-normal training set. Anomalies are then detected amongst unseen samples by measuring their closeness to the normal-data distribution.

Autoencoders are an extremely popular approach to learn the behaviour of normal data. The reconstruction error of a sample is directly used as its anomaly score; anomalies are assumed to have a higher reconstruction error than

normal samples due to their difference in distribution. However, this anomaly score fails to account for the fact that reconstruction error can vary greatly even amongst different types of normal samples. For example, consider a sensor system in a factory with various activities on weekdays but zero activity on weekends. The weekday samples are diverse and complex, therefore we could expect the reconstruction error of these samples to vary greatly. Meanwhile, even a small amount of activity in a weekend sample would be anomalous, even if the effect on the reconstruction error is minimal. The anomaly detector would likely detect false positives (high reconstruction error) among weekday samples and false negatives among weekend samples. Although encoded within the input data attributes, the context of each individual sample (i.e. day of the week) is neglected at the detection stage by the standard scoring approach. Instead, all samples are assessed according to the same error threshold or standard. This invokes the implicit assumption that the reconstruction errors of all samples are identically distributed, regardless of any individual characteristics and context which could potentially influence the reconstruction error significantly.

In this work, we aim to address this problem by proposing **Adaptive Reconstruction Error-based Scoring (ARES)**. Our **locally-adaptive** scoring method is able to automatically account for any contextual information which affects the reconstruction error, resulting in more accurate anomaly detection. We use a flexible, **neighbourhood-based** approach to define the context, based on the location of its latent representation learnt by the model.

Our scoring approach is applied at test time, so it can be retrofitted to pre-trained models of any size and architecture or used to complement existing anomaly detection techniques. Our score is simple and efficient to compute, requiring very little additional computational time, and the code is available online⁴. In summary, we make the following contributions:

1. We **empirically show with real, multi-class data** the variation in reconstruction error among different samples in the normal set and why this justifies the need for local adaptivity in anomaly scoring.
2. We propose a novel anomaly scoring method which adapts to this variation by evaluating anomaly status based on the local context of a given sample in the latent space.
3. We evaluate our method against a wide range of baselines with various benchmark datasets. We also study the effect of different components and formulations of our scoring method in an ablation study.

2 Background

2.1 Autoencoders

Autoencoders are neural networks that output a reconstruction of the input data [26]. They are comprised of two components: an encoder and decoder. The

⁴ <https://github.com/agoodge/ARES>

encoder compresses data from the input level into lower-dimensional latent representations through its hidden layers. The encoder output is typically known as the bottleneck, from which the reconstruction of the original data is found through the decoder hidden layers. The network is trained to minimise the reconstruction error over the training set:

$$\min_{\theta, \phi} \|\mathbf{x} - (f_{\theta} \circ g_{\phi})(\mathbf{x})\|_2^2, \tag{1}$$

where g_{ϕ} is the encoder, f_{θ} the decoder. The assumption is that data lies on a lower-dimensional manifold within the high-dimensional input space. The auto-encoder learns to reconstruct data on this manifold by performing dimension reduction. By training the model on only normal data, the reconstruction error of a normal sample should be low as it close to the learnt manifold on which it has been reconstructed by the autoencoder, whereas anomalies are far away and are reconstructed with higher error.

2.2 Local Outlier Factor

The Local Outlier Factor (LOF) method is a neighbourhood-based approach to anomaly detection; it measures the density of a given sample relative to its k nearest neighbours’ [8]. Anomalies are assumed to be in sparse regions far away from the one or more high-density clusters of normal data. A lower density therefore suggests the sample is anomalous. The original method uses the ‘reachability distance’; defined for a point A from point B as:

$$\max\{k\text{-distance}(B), d(A, B)\} \tag{2}$$

where $d(\cdot, \cdot)$ is a chosen distance metric and $k\text{-distance}(B)$ is the distance of B to its k^{th} nearest neighbour. The local reachability density and subsequently the local outlier factor of A based on its set of neighbours $\mathbf{N}_k(A)$, is:

$$\text{lrd}_k(A) := \left(\frac{\sum_{B \in \mathbf{N}_k(A)} \text{reachability-distance}_k(A, B)}{|\mathbf{N}_k(A)|} \right)^{-1}. \tag{3}$$

$$\text{LOF}_k(A) = \frac{\sum_{B \in \mathbf{N}_k(A)} \text{lrd}_k(B)}{|\mathbf{N}_k(A)| \cdot \text{lrd}_k(A)}, \tag{4}$$

3 Related Work

Anomalies are often assumed to occupy sparse regions far away from high-density clusters of normal data. As such, a great variety of methods detect anomalies by measuring the distances to their nearest neighbours, most notably KNN [35], which directly uses the distance of a point to its k^{th} nearest neighbour as the anomaly score. LOF [8] measures the density of a point relative to the density of points in its local neighbourhood, as seen in Section 2. This is beneficial as a

given density may be anomalous in one region but normal in another region. This local view accommodates the natural variation in density and therefore allows for the detection of more meaningful anomalies. More recently, deep models such as graph neural networks have been used to learn anomaly scores from neighbour-distances in a data-driven way [15]. Naturally, this relies on an effective distance metric. This is often itself a difficult problem; even the most established metrics have been shown to lose significance in high-dimensional spaces [6] due to the ‘curse of dimensionality’.

Reconstruction-based methods rely on deep models, such as autoencoders, to learn to reconstruct samples from the normal set accurately. Samples are then flagged as anomalous if their reconstruction error higher than some pre-defined threshold, as it is assumed that the model will reconstruct samples outside of the normal set with a higher reconstruction error [26,13,4,9,14].

More recent works have used autoencoders or other models are deep feature extractors, with the learnt features then used in another anomaly detection module downstream, such as Gaussian mixture models [7], DBSCAN [3], KNN [5] and auto-regressive models [1]. [11] calibrates different types of autoencoders against the effects of varying hardness within normal samples.

Other methods measure the distance of a sample to a normal set-enclosing hypersphere [30,25]. or its likelihood under a learnt model [12,24,34]. Generative adversarial networks have also been proposed for anomaly detection, mostly relying either on the discriminator score or the accuracy of the generator for a given sample to determine anomalousness [2,33].

In all of these methods, the normal set is often restricted to a subset of the available data; e.g. just one class from a multi-class dataset in experiments. In practice, normal data could belong to multiple classes or modes of behaviour which all need to be modelled adequately. As such, developing anomaly detection methods that can model a diverse range of normal behaviours and adapt their scoring appropriately are important.

4 Methodology

4.1 Problem Definition

We assume to have m normal training samples $\mathbf{x}_1^{(\text{train})}, \dots, \mathbf{x}_m^{(\text{train})} \in \mathbb{R}^d$ and n testing samples, $\mathbf{x}_1^{(\text{test})}, \dots, \mathbf{x}_n^{(\text{test})} \in \mathbb{R}^d$, each of which may be normal or anomalous. For each test sample \mathbf{x} , our algorithm should indicate how anomalous it is through computing its **anomaly score** $s(\mathbf{x})$. The goal is for anomalies to be given higher anomaly scores than normal points. In this work, the fundamental question is:

Given an autoencoder with encoder g_ϕ and decoder f_θ , how can we use the latent encoding $\mathbf{z} = g_\phi(\mathbf{x})$ and reconstruction $\hat{\mathbf{x}} = (f_\theta \circ g_\phi)(\mathbf{x})$ of a sample \mathbf{x} to score its anomalousness?

Our approach is not limited to standard autoencoders, and can be applied to other models that involve a latent encoding \mathbf{z} and a reconstruction $\hat{\mathbf{x}}$.

In practice, the anomaly score is compared with a user-defined anomaly threshold; samples which exceed this threshold are flagged as anomalies. Different approaches can be employed to set this threshold, such as extreme value theory [28]. Our focus is instead on the approach to anomaly scoring itself, which allows for any choice of thresholding scheme to be used alongside it.

4.2 Statistical Interpretation of Reconstruction-based Anomaly Detection

In the standard approach, with the residual as $\boldsymbol{\varepsilon} := \mathbf{x} - \hat{\mathbf{x}}$, the anomaly score is:

$$R(\boldsymbol{\varepsilon}) := \|\boldsymbol{\varepsilon}\|_2^2 = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2. \tag{5}$$

This can be seen as a negative log-likelihood-based score that assumes $\boldsymbol{\varepsilon}$ follows a Gaussian distribution with zero mean and unit variance: $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$:

$$-\log P(\boldsymbol{\varepsilon}) = \frac{d}{2} \log(2\pi) + \frac{R(\boldsymbol{\varepsilon})}{2} \tag{6}$$

The negative log-likelihood is an intuitive anomaly score because anomalous samples should have lower likelihood, thus higher negative log-likelihood.

The key conclusion, ignoring the constant additive and multiplicative factors, is that $R(\boldsymbol{\varepsilon})$ can be equivalently seen as a negative log-likelihood-based score, based on a model with the implicit assumption of constant mean as well as *constant variance residuals across all samples*, known in the statistical literature as **homoscedasticity**. Most crucially, this assumption applies to all \mathbf{x} regardless of its individual latent representation, \mathbf{z} , which implies that the reconstruction error of \mathbf{x} does not depend on the location of \mathbf{z} in the latent space. We question the validity of this assumption in real datasets.

In our earlier example of a factory sensor system, the variance of a sample \mathbf{x} depends greatly on its characteristics (i.e. day of the week), with high variance during weekdays and low variance on weekends. This is instead a prime example of **heteroscedasticity**. Furthermore, since the latent representation \mathbf{z} typically encodes these important characteristics of \mathbf{x} , there is likely to be a clear dependence between \mathbf{z} and the reconstruction error $R(\boldsymbol{\varepsilon})$, which will be examined in Section 4.3. This relationship could be exploited to improve the modelling of reconstruction errors and thus detection accuracy by adapting the anomaly score to this relationship at the sample-level.

4.3 Motivation & Empirical Results

In this section, we train an autoencoder on all 10 classes of MNIST to empirically examine the variation in the reconstruction error of real data to scrutinise the homoscedastic assumption. In doing so, we exemplify the shortcomings of the standard approach which implicitly assumes all reconstruction errors from a fixed Gaussian distribution. The autoencoder architecture and training procedure is the same as those described in Section 5. We make the following observations:

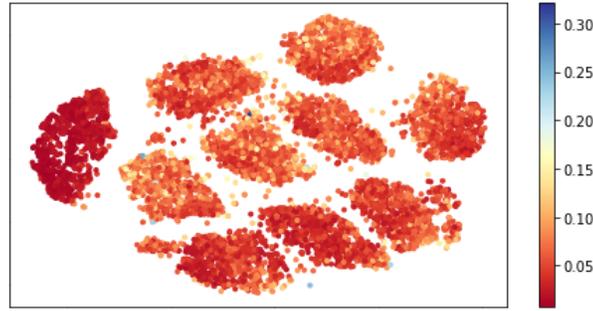


Fig. 1: t-SNE plot of the latent encodings with colour determined by reconstruction error of the associated sample.

1. **Inter-class variation:** In Figure 1, the t-SNE [21] projections of the latent representations of training points are coloured according to their reconstruction error. We see that the autoencoder learns to separate each class approximately into distinct clusters. In this context, the class label can be seen as a variable characteristic between different samples within the normal set (similar to the weekday vs. weekend example).

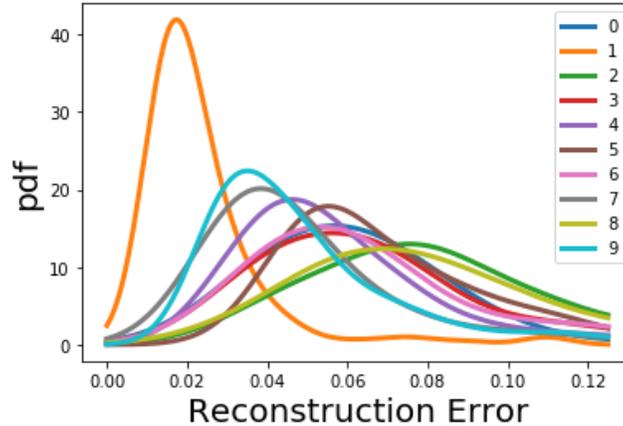


Fig. 2: Probability density function of reconstruction errors associated with different classes of training samples, estimated via kernel density estimation.

2. **Intra-class variation** Also in Figure 1, we can see that there is significant variation even within any single cluster. In other words, there are noticeably distinct regions of high (and low) reconstruction errors within each individual cluster. This shows that there are additional characteristics that influence

whether a given normal sample has a high or low reconstruction error besides its class label.

We observe that there is significant variation in reconstruction error between the different clusters. Most notably, samples in the leftmost cluster (corresponding to class 1) has significantly lower reconstruction errors than most others. Indeed, as shown in Figure 2, the distribution of reconstruction errors associated with class 1 samples is very different to those of the other classes. This shows that there is significant variation in reconstruction errors between classes, and that it is inappropriate to assume that the reconstruction errors of all samples can be modelled by a single, fixed Gaussian distribution. For example, a reconstruction error of 0.06 would be high for a class 1 or class 9 sample, but low for a class 2 or class 8 sample.

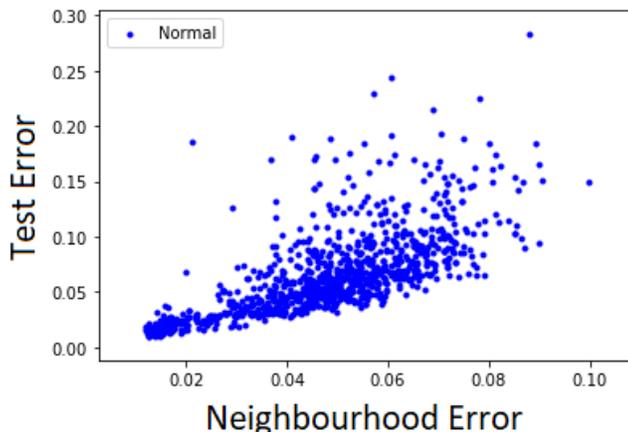


Fig. 3: Average reconstruction error of the training set nearest neighbours of a point versus its own reconstruction error for normal test samples.

3. **Neighbourhood correlation:** In Figure 3, we plot the reconstruction error of test samples against the average reconstruction error of its nearest training set neighbours in the latent space (neighbourhood error). We see that the reconstruction error of a test point increases as its neighbourhood error increases. Furthermore, the variance in test errors increases for larger neighbourhood errors: a clear heteroscedastic relationship. This information is useful in determining anomalousness: a test error of 0.1 would be anomalously high if its neighbourhood error is 0.02, but normal if it is 0.06.

Given these observations, we propose that anomalies could be more accurately detected by incorporating contextual information into the anomaly scoring beyond reconstruction error alone. Furthermore, analysing the neighbourhood of a given sample provides this contextual information to help determine its anoma-

lousness. In section 4.4, we propose our Adaptive Reconstruction Error-based Scoring method to achieve this.

4.4 Adaptive Reconstruction Error-based Scoring

In Section 4.2, we saw that the standard approach assumes that all residuals come from a fixed Gaussian with constant mean and unit variance: $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$. In Section 4.3, we saw that this assumption is inappropriate. In this section, we detail ARES, a novel anomaly scoring methodology which aims to address this flaw by adapting the scoring for each samples local context in the latent space.

The normal level of reconstruction error varies for samples in different regions of the latent space, meaning the latent encoding of a sample holds important information regarding anomalousness. As such, ARES is inspired by the **joint-likelihood** of a samples residual ε with its latent encoding \mathbf{z} , defined as follows:

$$-\log P(\varepsilon, \mathbf{z}) = -\log P(\varepsilon|\mathbf{z}) - \log P(\mathbf{z}). \quad (7)$$

The first term, $-\log P(\varepsilon|\mathbf{z})$, is the conditional (negative log-) likelihood of the points residual conditioned on its latent encoding. The second term, $-\log P(\mathbf{z})$, is the likelihood of observing the latent encoding from the normal set.

We now detail our approach to interpret these terms into tractable, efficient scores which we name the **local reconstruction score** and **local density score** respectively. These scores are combined to give the overall ARES anomaly score.

4.5 Local Reconstruction Score

The local reconstruction score is based on an estimate of how likely a given residual ε (and consequent reconstruction error) is to come from the corresponding sample \mathbf{x} , based on its latent encoding \mathbf{z} :

$$r(\mathbf{x}) = -\log P(\varepsilon|\mathbf{z}) \quad (8)$$

This likelihood cannot be calculated directly for any individual \mathbf{z} . Instead, we consider the k nearest neighbours of \mathbf{z} in the latent space, denoted $\mathbf{N}_k(\mathbf{z})$, to be a sample population which z belongs to. This is intuitive as data points with similar characteristics to \mathbf{x} in the input space are more likely to be encoded nearer to \mathbf{z} in the latent space. The full local reconstruction score algorithm is shown in Algorithm 1.

After training the autoencoder, we fix the model weights and store in memory all training samples' reconstruction errors and latent encodings taken from the bottleneck layer. For a test sample \mathbf{x} , we find its own latent encoding \mathbf{z} and reconstruction $\hat{\mathbf{x}}$. We then find $\mathbf{N}_k(\mathbf{z})$, the set of k nearest neighbours to \mathbf{z} amongst the training set encodings. We only find nearest neighbours among the training samples as they are all assumed to be normal and therefore should have reconstruction errors within the normal range.

We want to obtain an estimate of the reconstruction error that could be expected of \mathbf{x} , assuming its a normal point, based on the reconstruction errors of

its (normal) neighbours. We can then compare this expected value, conditioned on its unique location in the latent space, with the points true reconstruction error to determine its anomalousness.

In a fully probabilistic approach, we could do this by measuring the likelihood of the test points reconstruction error under a probability distribution, e.g. a Gaussian, fit to the neighbours’ reconstruction errors. However, it is unnecessarily restrictive to assume any closed-form probability distribution to adequately model this population for the unique neighbourhood of each and every test sample. Instead, we opt for a non-parametric approach; we measure the difference between the test points reconstruction error and the median reconstruction error of its neighbouring samples. The larger the difference between them, the more outlying the test point is in comparison to its local neighbours, therefore the more likely it is to be anomalous. In practice, using the median was found to perform better than the mean as it is more robust to extrema. With this, we obtain the local reconstruction score as:

$$r(\mathbf{x}) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 - \underset{\mathbf{n} \in \mathbf{N}_k(\mathbf{z})}{\text{median}}(\|\mathbf{n} - \hat{\mathbf{n}}\|_2^2). \tag{9}$$

Nearest neighbour search in the latent space is preferable over the input space as dimension reduction helps to alleviate the curse of dimensionality (see Section 3), resulting in more semantically-meaningful neighbors. Secondly, from a practical perspective, the neighbour search is less time-consuming and computationally intensive in lower dimensional spaces.

Algorithm 1 Local Reconstruction Score

Input: Autoencoder $A(\cdot) = (f_\theta \circ g_\phi)(\cdot)$, training set $\mathbf{X}_{\text{train}}$, test sample \mathbf{x}_{test}

Parameters: neighbour count k

Output: Local reconstruction score $r(\mathbf{x}_{\text{test}})$

- 1: Train autoencoder A on training set according to: $\min_{\theta, \phi} \|\mathbf{X}_{\text{train}} - \hat{\mathbf{X}}_{\text{train}}\|_2^2$ where $\hat{\mathbf{x}}_{\text{train}}^i = (f_\theta \circ g_\phi)(\mathbf{x}_{\text{train}}^i)$ for $\mathbf{x}_{\text{train}}^i \in \mathbf{X}_{\text{train}}$.
 - 2: Find latent encoding of \mathbf{x}_{test} : $\mathbf{z}_{\text{test}} := g_\phi(\mathbf{x}_{\text{test}})$
 - 3: Find the set of k nearest neighbours to \mathbf{z}_{test} among latent encodings of the training data: $\mathbf{N}_k(\mathbf{z}_{\text{test}}) := \{\mathbf{z}_{\text{train}}^1, \dots, \mathbf{z}_{\text{train}}^k\}$ where $\mathbf{z}_{\text{train}}^i = g_\phi(\mathbf{x}_{\text{train}}^i) \in$ for $i = \{1, \dots, k\}$.
 - 4: **return** $r(\mathbf{x}_{\text{test}}) = \|\mathbf{x}_{\text{test}} - \hat{\mathbf{x}}_{\text{test}}\|_2^2 - \underset{\mathbf{n} \in \mathbf{N}_k(\mathbf{z})}{\text{median}}(\|\mathbf{x}_{\text{train}}^i - \hat{\mathbf{x}}_{\text{train}}^i\|_2^2)$
-

4.6 Local Density Score

The local reconstruction score corresponds to the conditional term of the joint likelihood. We now introduce the local density score, which corresponds to the likelihood of observing the given encoding in the latent space. This is a density estimation task, which concerns the relative distance of \mathbf{z} to its nearest neighbours, unlike the local reconstruction score which focuses on the reconstruction

error of neighbours. Anomalies are assumed to exist in sparse regions, where normal samples are unlikely to be found in significant numbers. Any multivariate distribution P , with trainable parameters Θ , could be used to estimate this density:

$$d(\mathbf{x}) := -\log P(\mathbf{z}; \Theta) \quad (10)$$

Note that it is common to ignore constant factor shifts in the anomaly score. Thus, the distribution P need not be normalized; even unnormalized density estimation techniques can be used as scoring functions. We note that LOF is an example of an unnormalized score which is similarly locally adaptive like the local reconstruction score, so LOF is used for the local density score in our main experiments. Other methods are also tested and their performance is shown in the ablation study.

The overall anomaly score for sample \mathbf{x} is:

$$s(\mathbf{x}) := r(\mathbf{x}) + \alpha d(\mathbf{x}), \quad (11)$$

where $r(\mathbf{x})$ is its local reconstruction score and $d(\mathbf{x})$ the local density score. These two scores are unnormalized, so we use a scaling factor α to balance the relative magnitudes of the two scores. We heuristically set it equal to 0.5 for all datasets and settings for simplicity, as this was found to balance the two scores sufficiently fairly in most cases. We choose not to treat α as a hyper-parameter to be tuned to optimise performance, although different values could give better performance for different datasets. The effect of changing α is shown in the supplementary material.

Computational Runtime: The average runtimes of experiments with the MNIST dataset can be found in the supplementary material. We see that, despite taking longer than the standard reconstruction error approach, the additional computational runtime of ARES anomaly scoring is insignificant in relation to the model training time. Anomaly scoring with ARES is just 1.2% of the overall time taken (including training) in the one-class case (0.017 minutes for scoring versus 1.474 minutes for training), and 3.5% in the multi-class case. The additional run-time is a result of the k nearest neighbour search. An exact search algorithms would be $\mathcal{O}(nm)$ with n train and m test samples, however approximate methods can achieve near-exact accuracy much more efficiently.

5 Experiments

In our experiments, we aim to answer the following research questions:

RQ1 (Accuracy): Does ARES perform better than existing anomaly detection methods?

RQ2 (Ablation Study): How do different components and design choices of ARES contribute to its performance?

5.1 Datasets and Experimental Setup

Dataset	#Dim	#Classes	#Samples	Description
SNSR [31]	48	Multi-class	58,509	Electric current signals
MNIST [19]	784	Multi-class	70,000	0-9 digit images
FMNIST [32]	784	Multi-class	70,000	Fashion article images
OTTO [22]	93	Multi-class	61,878	E-commerce types
MI-F [29]	58	Single-class	25,286	CNC milling defects
MI-V [29]	58	Single-class	23,125	CNC milling defects
EOPT [16]	20	Single-class	90,515	Storage system failures

Table 1: Name and descriptions of the datasets used in experiments, including the number of samples and dimensions.

Table 1 shows the datasets we use in experiments. The single-class datasets consist of ground truth normal-vs-anomaly labels, as opposed to the multiple class labels in multi-class datasets. In the latter, we distinguish between the ‘one-class normality’ setup, in which one class label is used as the normal class and all other classes are anomalous. Alternatively, in the ‘multi-class normality’ setup, one class is anomalous and all other classes are normal. For a dataset with N classes, there are N possible arrangements of normal and anomaly classes. We train separate models for each arrangement and find their average score for the final result. We use the Area-Under-Curve (AUC) metric to measure performance as it does not require an anomaly score threshold to be set.

As anomaly scores are calculated for each test sample independently of each other, the proportion of anomalies in the test set has no impact on the anomaly score assigned to any given sample. Therefore, we are able to use a normal:anomaly ratio of 50:50 in our experiments for the sake of simplicity and an unbiased AUC metric. Besides the normal samples in the test set, the remaining normal samples are split 80:20 into training and validation sets for all models. Full implementation details can be found in the supplementary material.

5.2 Baselines

We test the performance of ARES against a range of baselines. We use the scikit learn implementations of LOF (in the input space) [8], IFOREST [20], PCA [27] and OC-SVM [10]. Publicly available codes are used for DAGMM [7], RAPP-SAP and RAPP-NAP [17], and we use Pytorch to build the autoencoder (AE and ARES) and variational autoencoder (VAE). All experiments were conducted in Windows OS using an Nvidia GeForce RTX 2080 Ti GPU.

We do not tune hyper-parameters relating to the model architectures or training procedures for any method. The effect of variation in hyper-parameters is studied in the ablation study in Section 5.4 and the supplementary material

instead. We set the number of neighbours $k = 10$ for both the local density and local reconstruction score in our main experiments.

5.3 RQ1 (Accuracy):

Table 2 shows the average AUC scores as a percentage (i.e. multiplied by 100). In the one-class normality setting, ARES significantly improves performance over the baselines in all multi-class datasets besides FMNIST. In the single-class datasets, this improvement is even greater, e.g. +8% lift for MI-F and EOPT. Compared with AE, we see that local adaptivity helps to detect true anomalies by correcting for the natural variation in reconstruction error in the latent space. This effect is even more pronounced in the multi-class normality setting, where ARES gives the best performance on all datasets. Here, the normal set is much more diverse and therefore local adaptivity is even more important. We show the standard deviations and additional significance test scores in the supplementary material, which also show statistically significant ($p < 0.01$) improvement.

Dataset	LOF	IFOREST	PCA	OC-SVM	SAP	NAP	DAGMM	VAE	AE	ARES
One-class Normality										
SNSR	97.98	89.16	92.01	95.85	98.79	98.74	88.08	89.49	98.30	98.83**
MNIST	96.85	85.44	95.68	90.35	95.35	97.25	89.60	91.73	96.96	97.89**
FMNIST	91.35	91.39	90.13	90.74	89.66	93.08	87.97	77.59	92.33	91.63
OTTO	84.76	70.34	80.09	81.43	81.61	82.77	68.02	82.39	85.26	87.86**
MI-F	59.79	81.53	55.07	76.69	81.78	80.61	82.23	76.93	71.19	89.52
MI-V	83.97	84.35	87.32	83.58	88.24	89.35	75.45	89.03	90.75	93.94
EOPT	55.01	61.61	54.72	59.66	59.87	61.69	60.63	68.08	59.85	68.43
Multi-class Normality										
SNSR	60.74	52.70	52.94	52.79	57.52	58.32	54.77	61.36	57.28	69.78**
MNIST	77.40	56.49	70.41	58.56	84.76	86.38	54.24	84.12	80.04	93.25**
FMNIST	71.50	64.75	66.92	60.27	68.50	72.09	57.56	71.18	71.03	72.49
OTTO	63.01	54.14	58.27	62.96	57.47	63.44	58.96	61.88	59.59	63.54

Table 2: Mean AUC scores for each datasets and normality setting. The best scores are highlighted in bold and we mark the most significant improvements over AE ($p < 0.01$) with **. Further tests are in the supplementary material.

The neighbours of a normal sample with high reconstruction error tend to be have high reconstruction errors themselves. By basing the anomaly scoring on their relative difference, ARES uses this to better detect truly anomalous samples. Furthermore, ARES also uses the local density of the point, which depends purely on its distance to the training samples in the latent space. This is important as there may be some anomalies with such low reconstruction error that comparison with neighbours does not alone indicate anomalousness (for example

the weekend samples mentioned earlier). These samples could be expected to be occupy very sparse regions of the latent space due to their significant deviation from the normal set, which means they can be better detected through the local density score. By combining these two scores, we are able to detect a wider range of anomalous data than either could individually.

5.4 RQ2 (Ablation Study):

Density Estimation Method Table 3 shows the performance of ARES with other density estimation methods. KNN is the distance to the k^{th} nearest neighbour in the latent space with $k = 20$. GD is the distance of a point to the closest of N Gaussian distributions fit to the latent encodings of samples for each of the classes in the training set ($N = 1$ in the one-class normality case). NF is the likelihood under a RealNVP normalizing flow [23].

Dataset	LOF	KNN	GD	NF
One-class Normality				
SNSR	98.83	98.66	95.68	98.29
MNIST	97.89	95.24	87.29	97.30
FMNIST	91.63	92.94	90.26	92.14
OTTO	87.76	83.14	81.27	86.74
MI-F	89.52	76.12	80.41	84.47
MI-V	93.94	91.60	79.61	92.89
EOPT	68.43	67.63	63.35	61.07
Multi-class Normality				
SNSR	69.78	67.42	61.86	60.63
MNIST	93.25	92.92	91.14	86.25
FMNIST	72.49	72.77	73.02	70.48
OTTO	63.54	61.71	66.30	63.67

Table 3: Mean AUC scores for each choice of local density score. The best scores are highlighted in bold.

LOF performs best overall, closely followed by KNN. GD performs poorly in one-class experiments, however it is better in multi-class normality experiments and even the best for FMNIST and OTTO. This could be as the use of multiple distributions provides more flexibility. NF is noticeably worse; previous studies have found that normalizing flows are not well-suited to detect out-of-distribution data [18].

In the supplementary material, we further test these density estimation methods by varying their hyper-parameters and find LOF to still come out best. Further supplementary experiments show that the local density score generally performs better than the local reconstruction score in the multi-class normality

case and vice versa in the one-class case. Combining them, as in ARES, generally gives the best performance overall across different latent embedding sizes.

Robustness to Training Contamination In practice, it is likely that a small proportion of anomalies ‘contaminate’ the training set. In this section, we study the effect of different levels of contamination on ARES, defined as $n\%$ of the total number of samples in the training set. Table 4 shows that increasing n worsens performance overall. With more anomalies in the training set, it is more likely that anomalies are found in the nearest neighbour set of more test points, which skews both the average neighbourhood reconstruction error as well as their density in the latent space and degrades performance.

Neighbourhood Size					
n	10	50	100	200	500
One-class Normality					
0	97.89	97.85	97.80	97.73	97.56
0.5	95.77	96.08	96.06	95.95	95.72
1	82.50	94.45	95.97	95.81	95.42
2	90.65	92.88	93.29	93.33	92.99
3	88.41	91.35	91.92	92.07	91.86
5	84.91	87.85	89.20	90.26	90.65
10	79.69	81.51	83.16	84.88	86.82
Multi-class Normality					
0	93.25	92.71	92.08	91.31	89.68
0.5	68.38	76.07	78.36	80.75	81.42
1	52.38	58.42	67.62	65.75	68.91
2	59.50	63.13	65.71	67.98	69.68
3	57.09	59.00	60.78	62.89	65.35
5	50.86	51.48	52.40	53.47	55.78
10	52.50	50.80	50.87	51.14	52.11

Table 4: Mean AUC scores in MNIST with training set anomaly contamination ($n\%$) and different neighbourhood sizes. The best scores are highlighted in bold.

We find that ARES is more robust to training set contamination with a higher setting of the number of neighbours (k). In the one-class setup, we see that ARES performs better with higher values of k as the proportion of anomalies increases. By using more neighbours, the effect of any individual anomalies on the overall neighbourhood error is reduced, which helps to maintain better performance. This effect is even more stronger in the multi-class normality setup. The highest neighbour count of $k = 500$ gives the best performance in all cases except for $n = 0\%$ and 10% . As the multi-class training sets are much larger than their

one-class normality counterparts, $n\%$ corresponds to a much larger number of anomalous contaminants, which explains their greater effect for a given k .

6 Conclusion

Autoencoders are extremely popular deep learning models used for anomaly detection through their reconstruction error. We have shown that the assumption made by the standard reconstruction error score, that reconstruction errors are identically distributed for all normal samples, is unsuitable for real datasets. We empirically show that there is a heteroscedastic relationship between latent space characteristics and reconstruction error, which demonstrates why adaptivity to local latent information is important for anomaly scoring. As such, we have developed a novel approach to anomaly scoring which adaptively evaluates the anomalousness of a samples reconstruction error, as well as its density in the latent space, relative to those of its nearest neighbours. We show that our approach results in significant performance improvements over the standard approach, as well as other prominent baselines, across a range of real datasets.

Acknowledgements

This work was supported in part by NUS ODPRT Grant R252-000-A81-133.

References

1. Abati, D., Porrello, A., Calderara, S., Cucchiara, R.: Latent space autoregression for novelty detection. In: ICCV. pp. 481–490 (2019)
2. Akcay, S., Atapour-Abarghouei, A., Breckon, T.P.: Ganomaly: Semi-supervised anomaly detection via adversarial training. In: ACCV. pp. 622–637. Springer (2018)
3. Amarbayasgalan, T., Jargalsaikhan, B., Ryu, K.H.: Unsupervised novelty detection using deep autoencoders with density based clustering. Applied Sciences **8**(9), 1468 (2018)
4. An, J.: Variational Autoencoder based Anomaly Detection using Reconstruction Probability. In: SNU Data Mining Center 2015-2 Special Lecture on IE (2015)
5. Bergman, L., Cohen, N., Hoshen, Y.: Deep nearest neighbor anomaly detection. arXiv preprint arXiv:2002.10445 (2020)
6. Beyer, K., Goldstein, J., Ramakrishnan, R., Shaft, U.: When is “nearest neighbor” meaningful? In: International conference on database theory. pp. 217–235. Springer (1999)
7. Bo, Z., Song, Q., Chen, H.: Deep autoencoding gaussian mixture model for unsupervised anomaly detection (2018)
8. Breunig, M.M., Kriegel, H.P., Ng, R.T., Sander, J.: LOF: identifying density-based local outliers. In: SIGMOD. vol. 29, pp. 93–104. ACM (2000)
9. Chen, J., Sathe, S., Aggarwal, C., Turaga, D.: Outlier detection with autoencoder ensembles. In: SDM. pp. 90–98. SIAM (2017)
10. Chen, Y., Zhou, X.S., Huang, T.S.: One-class SVM for learning in image retrieval. In: ICIP. pp. 34–37. Citeseer (2001)

11. Deng, A., Goodge, A., Lang, Y.A., Hooi, B.: CADET: Calibrated Anomaly Detection for Mitigating Hardness Bias. *IJCAI* (2022)
12. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real nvp. *arXiv preprint arXiv:1605.08803* (2016)
13. Feng, W., Han, C.: A Novel Approach for Trajectory Feature Representation and Anomalous Trajectory Detection. *ISIF* pp. 1093–1099 (2015)
14. Goodge, A., Hooi, B., Ng, S.K., Ng, W.S.: Robustness of Autoencoders for Anomaly Detection Under Adversarial Impact. *IJCAI* (2020)
15. Goodge, A., Hooi, B., Ng, S.K., Ng, W.S.: Lunar: Unifying local outlier detection methods via graph neural networks. In: *Proceedings of the AAAI Conference on Artificial Intelligence* (2022)
16. inIT: Tool wear detection in cnc mill (2018), <https://www.kaggle.com/init-owl/high-storage-system-data-for-energy-optimization>
17. Kim, K.H., Shim, S., Lim, Y., Jeon, J., Choi, J., Kim, B., Yoon, A.S.: RaPP: Novelty Detection with Reconstruction along Projection Pathway. In: *ICLR* (2019)
18. Kirichenko, P., Izmailov, P., Wilson, A.G.: Why normalizing flows fail to detect out-of-distribution data. *arXiv preprint arXiv:2006.08545* (2020)
19. Lecun, Y.: Mnist (2012), <http://yann.lecun.com/exdb/mnist/>
20. Liu, F.T., Ting, K.M., Zhou, Z.H.: Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* **6**(1), 1–39 (2012)
21. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(Nov), 2579–2605 (2008)
22. Otto, G.: Otto group product classification challenge (2015), <https://www.kaggle.com/c/otto-group-product-classification-challenge>
23. Papamakarios, G., Pavlakou, T., Murray, I.: Masked autoregressive flow for density estimation. In: *NeurIPS*. pp. 2338–2347 (2017)
24. Rezende, D.J., Mohamed, S.: Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770* (2015)
25. Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Siddiqui, S.A., Binder, A., Müller, E., Kloft, M.: Deep one-class classification. In: *ICML*. pp. 4393–4402 (2018)
26. Sakurada, M., Yairi, T.: Anomaly detection using autoencoders with nonlinear dimensionality reduction. In: *MLSDA*. p. 4. *ACM* (2014)
27. Shyu, M.L., Chen, S.C., Sarinapakorn, K., Chang, L.: A novel anomaly detection scheme based on principal component classifier. *Tech. rep.*, Miami Univ Coral Gables FL Dept of Electric and Computer Engineering (2003)
28. Siffer, A., Fouque, P.A., Termier, A., Largouet, C.: Anomaly detection in streams with extreme value theory. In: *SIGKDD*. pp. 1067–1075 (2017)
29. SMART: Tool wear detection in cnc mill (2018), <https://www.kaggle.com/shasun/tool-wear-detection-in-cnc-mill>
30. Tax, D.M., Duin, R.P.: Support vector data description. *Machine learning* **54**(1), 45–66 (2004)
31. UCI: Sensorless drive diagnosis (2015)
32. Xiao, H., Rasul, K., Vollgraf, R.: Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms (2017)
33. Zenati, H., Foo, C.S., Lecouat, B., Manek, G., Chandrasekha, V.R.: Efficient gan-based anomaly detection (2019)
34. Zhai, S., Cheng, Y., Lu, W., Zhang, Z.: Deep structured energy based models for anomaly detection. *ICML* **48**, 1100–1109 (2016)
35. Zimek, A., Gaudet, M., Campello, R.J.G.B., Sander, J.: Subsampling for efficient and effective unsupervised outlier detection ensembles. In: *SIGKDD*. pp. 428–436 (2013)

Supplementary Material: Appendix

Experiments

Model Architecture and Training

The autoencoder models for the two image datasets had seven hidden layers in both the encoder and decoder, for 14 layers total, while the tabular datasets had six in each (12 total). All linear layers were followed by a Leaky ReLU layer and Batch Normalization except for the final output layer in the decoder. For the image datasets, the hidden layer sizes for the encoder were as follows: [784, 600, 500, 400, 300, 200, 100, 20] and the same in reverse order for the decoder. For tabular datasets, the bottleneck size was set equal to the number of PCA principal components needed to explain 90% of the data, and the hidden layer sizes decrease (increase) linearly in the encoder (decoder) accordingly. All autoencoder models (including VAE) models followed this structure for the sake of consistency. All models were trained for 350 epochs with early stopping activated if the loss function, mean squared error, did not achieve a new minimum for 20 consecutive epochs. The batch size was set at 250 samples.

Standard Deviation in AUC Score

Table 1 shows the standard deviations of the AUC scores over the N trials for each dataset and settings, where N is the number of classes in the dataset. Naturally, as each setup has a different set of normal and anomalous classes, the scores vary greatly. Therefore we show the ificance values for performance improvement in the next section.

One-sided Wilcoxon Significance Test

In Table 2, we show the p -values of the one-sided Wilcoxon signed-rank test [1]. We test our method against both the AE, to see the effect of local adaptivity, as well as RAPP-NAP performance, as it was generally the next best-performing method. A lower p -value indicates more confidence that our method performs better over the corresponding baseline method. The binary-class datasets are not shown as they only had one experimental setup, i.e. one AUC score to measure. The p -values for the SNSR, MNIST and OTTO datasets in particular are extremely low. The performance is more comparable between ARES and the baselines for FMNIST, hence the larger p -values. The median p -value of ARES against AE is 0.0073 and for RAPP-NAP it is 0.016.

Dataset	PCA	LOF	OC-SVM	AE	SAP	NAP	DAGMM	ARES
Multi-Class Normality								
SNSR	7.91	18.47	19.16	17.25	18.52	20.19	13.15	16.73
MNIST	16.48	14.91	16.65	22.64	13.57	12.54	14.30	7.77
FMNIST	16.12	17.01	17.10	17.43	17.48	15.20	23.79	16.23
OTTO	15.24	22.76	19.58	17.62	17.88	20.32	15.64	18.09
One-Class Normality								
SNSR	5.41	1.12	2.54	1.32	0.66	0.67	6.38	0.94
MNIST	4.11	2.72	6.59	2.44	2.82	1.90	5.76	1.75
FMNIST	6.25	4.77	5.74	5.07	6.12	4.84	8.62	4.04
OTTO	8.64	6.46	7.17	6.15	7.63	7.27	8.44	4.98

Table 1: Standard deviation of AUC scores over the different runs.

Table 2: p values from the one-sided Wilcoxon signed-rank test. Values lower than 0.05 and 0.01 are marked with * and ** respectively.

<i>p</i> value	AE	RAPP-NAP
One-Class Normality		
SNSR	0.0063**	0.2969
MNIST	0.0083**	0.0234*
FMNIST	0.9303	0.9937
OTTO	0.0038**	0.0038**
Multi-Class Normality		
SNSR	0.0017**	0.0029**
MNIST	0.0035**	0.0035**
FMNIST	0.2538	0.6006
OTTO	0.0693	0.6165

Runtime

Table 3 shows the runtimes of the autoencoder training, the regular AE anomaly scoring and our ARES scoring. We see that, despite noticeably longer computational time for our method at test time, it is very insignificant in relation to the overall pipeline of the model including training time.

Time	Training	AE score	ARES score
One-Class	1.47444	0.00014	0.01721
Multi-Class	24.29161	0.00037	0.87681

Table 3: Runtimes (minutes) of training and testing for AE and ARES for comparison for each normality setting. Runtimes are calculated for the MNIST dataset and averaged over each of the trials for each modality setup.

Analysis

To understand why local adaptivity improves anomaly detection performance, we study one particular problem setup: MNIST samples of class 8 are normal and all other classes are all anomalies. Figure 1 shows the reconstruction error of the test samples of all classes in this setup. We see that reconstruction errors of the class 8 samples are generally lower than those from other classes except for class 1. This means that, based on reconstruction error alone, many anomalies from class 1 will be falsely classified as normal.

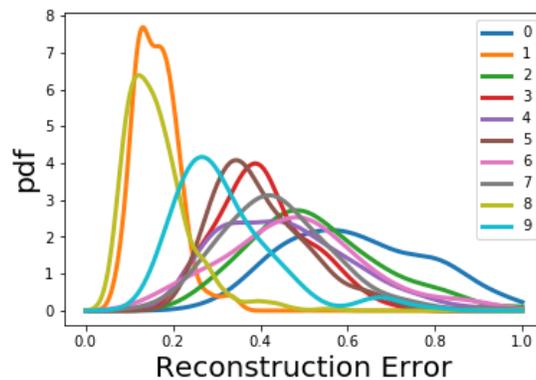


Fig. 1: Probability density functions of reconstruction errors of each class (class 8 is normal and others are anomalous).

ARES compares a points reconstruction error against its nearest neighbours which provides local context for more accurate detection. As shown in Figure 3, the neighbours of a poorly reconstructed normal test point tend to be the more poorly reconstructed training points, perhaps because these training points also contained some kind of abnormality. With this, the abnormality of these normal test points is less outlying within the context of the training points around them. This context is neglected in the standard scoring approach, and so is likely to detect false positives for these kinds of samples, unlike in our method which accounts for this context.

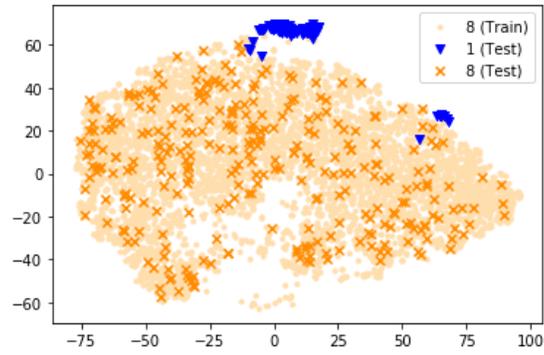


Fig. 2: t-SNE embeddings of class 1 test samples (blue), class 8 test samples (dark orange) and class 8 training samples (light orange).

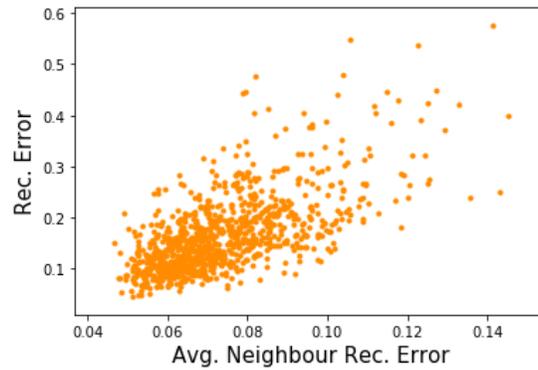


Fig. 3: Reconstruction error of normal (8) test samples versus the average reconstruction error of their training set nearest neighbours.

Furthermore, ARES also uses the local density of the point, which depends not on its reconstruction error but purely on its distance to the training samples in the latent space. In Figure 2, we can see that the vast majority of anomalous class 1 samples (in blue) are far away from their nearest neighbours in the training set than the normal class 8 test samples (in dark orange). By combining these two scores, we are able to detect such anomalies even if the reconstruction error alone cannot.

Scaling factor In Table 4, we see the average AUC score found when using different values for the scaling factor α . A lower or higher value of the scaling factor gives more weight to the local reconstruction or density score respectively.

Generally, we see that a lower value of α is better in the one-class normality case, while higher is better for the multi-class normality case, meaning the local density score is more beneficial when the diversity in the normal set is larger. This corroborates the results seen in Table 5, as the local reconstruction score is relatively better than the local density score in the one-class normality case and vice versa in the multi-class normality case.

α	One-Class Normality		Multi-Class Normality	
	MNIST	FMNIST	MNIST	FMNIST
0.1	97.51	91.52	88.08	70.84
0.25	97.76	91.70	91.53	71.97
0.5	97.89	91.63	93.25	72.49
1	97.72	91.16	93.91	72.49
2	97.00	90.25	93.90	72.17

Table 4: Average AUC score (%) of ARES with different settings of scaling factor α

Score Comparison Table 5 shows the performance of the local density and local reconstruction scores separately. We see that the local reconstruction score tends to do better in the one-class normality experiments while the local density score is better for multi-class normality experiments. However, combining both of them as in ARES gives the best performance in most cases overall.

Table 5: AUC score of the two components of ARES in isolation.

Dataset	$d(\mathbf{x})$	$r(\mathbf{x})$	ARES
One-Class Normality			
SNSR	96.24	98.49	98.83
MNIST	91.59	97.21	97.89
FMNIST	86.97	91.17	91.63
OTTO	77.17	88.57	87.86
MI-F	87.08	68.69	89.52
MI-V	87.93	94.73	93.94
EOPT	67.88	60.13	68.43
Multi-Class Normality			
SNSR	70.46	62.66	69.78
MNIST	93.23	81.96	93.25
FMNIST	71.05	68.97	72.49
OTTO	62.54	60.84	63.54

Density Estimation Method Table 6 shows the performance of ARES with other density estimation methods. KNN is the distance to the k^{th} nearest neighbour in the latent space. GD is the distance of a point to a Gaussian distribution fit to the training set latent encodings or the closest of N in the multimodal case. GD performs poorly in unimodal experiments, however it is better in multimodal experiments and even the best for FMNIST and OTTO. This could be as there are N separate distributions used for each of the N normal classes in this case. Within GD, we test the use of the Mahalanobis distance versus the Euclidean distance, and we see the former outperform the latter in most cases. NF is the likelihood under a RealNVP normalizing flow [2]. Other flows were tested but did not give stable results. It is noticeably worse than the other methods; previous studies have found that normalizing flows are not well-suited to detect out-of-distribution data.

Dataset		SNSR	MNIST	FMNIST	OTTO	MI-F	MI-V	EOPT
One-Class Normality								
LOF	$k = 10$	98.83	97.89	91.63	87.76	89.52	93.94	68.43
	$k = 40$	98.75	97.85	92.53	88.04	79.31	91.99	62.28
	$k = 100$	98.55	97.80	93.02	88.08	69.69	89.98	60.80
KNN	$k = 5$	98.65	95.53	92.47	83.99	77.26	93.55	79.45
	$k = 20$	98.66	95.24	92.94	83.14	76.12	91.60	67.63
	$k = 40$	98.56	94.93	92.99	82.76	74.24	90.508	64.91
GD	Euclidean	97.21	72.74	79.47	82.08	80.80	75.77	61.11
	Mahalanobis	95.68	87.29	90.26	81.27	80.41	79.61	63.35
NF		98.29	97.30	92.14	86.74	84.47	92.89	61.07
Multi-Class Normality								
LOF	$k = 10$	69.78	93.25	72.49	63.54			
	$k = 40$	67.65	92.90	72.74	63.80			
	$k = 100$	66.36	92.08	72.48	64.73			
KNN	$k = 5$	68.63	93.68	73.52	61.76			
	$k = 20$	67.42	92.92	72.77	61.74			
	$k = 40$	66.50	92.06	72.28	61.74			
GD	Euclidean	56.55	74.40	66.00	59.38			
	Mahalanobis	61.86	91.14	73.02	66.30			
NF		60.63	86.25	70.48	63.67			

Table 6: Mean AUC scores for each choice of local density score. The best scores are highlighted in bold.

Embedding Size We also study the effect of latent size, or the dimensionality of latent encodings, on the performance of our approach as well as the following variants:

- **AE**: The standard AE method without local adaptivity.
- **ARES-G**: AE reconstruction score combined with local density score.
- **L-R**: Local reconstruction score only.
- **L-D**: Local density score only.
- **ARES**: Combined local reconstruction and density score.

In Figure 4, we see this effect for the MNIST dataset in particular. We see that embedding size has little effect on reconstruction scores in the one-class normality setting but reduces performances significantly in the multi-class normality setting. In both settings, however, the density scores mostly improve with increasing embedding size. Overall, we see that combining the local reconstruction score with the density score gives the best performance across the board; it is more robust to variation in hyper-parameter choices.

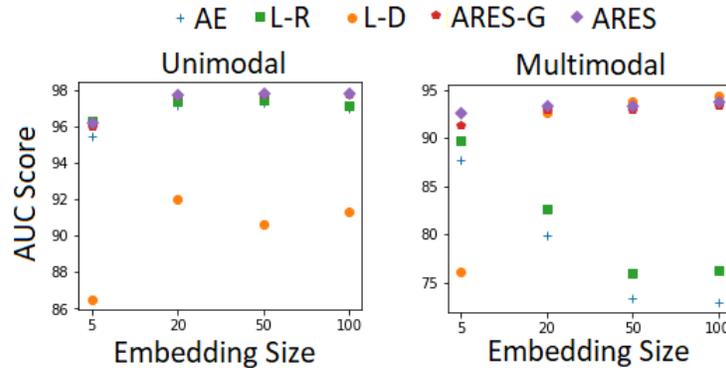


Fig. 4: Performance of various component methods across a range of embedding sizes.

References

1. Demsar, J.: Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research* 7, 1–30 (2006)
2. Papamakarios, G., Pavlakou, T., Murray, I.: Masked autoregressive flow for density estimation. In: *NeurIPS*. pp. 2338–2347 (2017)