
Normality-Calibrated Autoencoder for Unsupervised Anomaly Detection on Data Contamination

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we propose Normality-Calibrated Autoencoder (NCAE), which
2 can boost anomaly detection performance on the contaminated datasets with-
3 out any prior information or explicit abnormal samples in the training phase.
4 The NCAE adversarially generates high confident normal samples from a la-
5 tent space having low entropy and leverages them to predict abnormal samples
6 in a training dataset. NCAE is trained to minimise reconstruction errors in un-
7 contaminated samples and maximise reconstruction errors in contaminated sam-
8 ples. The experimental results demonstrate that our method outperforms shal-
9 low, hybrid, and deep methods for unsupervised anomaly detection and achieves
10 comparable performance compared with semi-supervised methods using labelled
11 anomaly samples in the training phase. The source code is publicly available on
12 https://github.com/nonamescientist/NCAE_UAD.git.

13 1 Introduction

14 Most of anomaly detection (AD) methods [Erfani et al., 2016, Zhai et al., 2016, Chen et al., 2017,
15 Ruff et al., 2018, Deecke et al., 2018, Ruff et al., 2019b, Golan and El-Yaniv, 2018, Pang et al.,
16 2019, Hendrycks et al., 2019a,b, Zong et al., 2018] assume that the training dataset only consists of
17 normal samples; however, datasets in real-world are easily *contaminated*, which means that datasets
18 contains both normal and abnormal samples. The contaminated samples significantly degrade the
19 AD performance of models derived based on the assumption.

20 Various methods have been proposed [Ruff et al., 2019a, Song et al., 2017, Akcay et al., 2018,
21 Chalapathy and Chawla, 2019, Zong et al., 2018] to improve the robustness of AD methods on
22 contaminated datasets. Particularly, filtering contaminated samples based on contamination ratio
23 [Zong et al., 2018, Ruff et al., 2019a], semi-supervised learning approaches that uses explicit abnormal
24 samples in the training step [Wang et al., 2005, Liu and Zheng, 2006, Görnitz et al., 2013, Ruff et al.,
25 2019a], and contamination sample prediction approaches based on geometric distance measurement
26 [Berg et al., 2019, Li et al., 2021, Lai et al., 2020], have been proposed. However, the aforementioned
27 approaches are domain or data-type specific. Additionally, those methods assume that abnormal
28 samples are likely to be located far from the distribution of normal samples, and the entropy of
29 abnormal samples is higher than that of normal samples [Berg et al., 2019, Li et al., 2021, Lai
30 et al., 2020]. Unfortunately, as shown in Figure 1, if a training dataset is highly contaminated, the
31 contaminated samples can also form a low entropy space by themselves.

32 In this paper, we present Normality-Calibrated Autoencoder (NCAE), which is robust to the training
33 dataset contamination. Our key idea on the NCAE is to adversarially generate high confident normal
34 samples from a low entropy feature space and then to contrastively compare the generated samples
35 with the input samples for estimating contamination score. After identifying the contaminated
36 samples, NCAE is trained to maximise reconstruction error of the found sample.

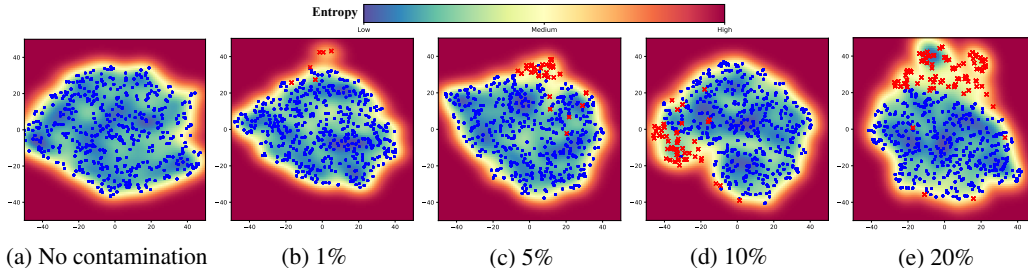


Figure 1: Entropy and distribution of latent features under different contamination ratios: (a) No contamination, (b) 1%, (c) 5%, (d) 10%, and (e) 20%. The samples on the ‘5’ class on the MNIST dataset are used as normal (blue dots) and contaminated samples (red x-marks) are randomly picked from the training samples of the remaining classes. The 500 samples are randomly picked for the visualisation. When a dataset is highly contaminated (*i.e.*, contamination ratio over 10%), contaminated samples are also located in a low entropy region.

2 Normality-Calibrated Autoencoder

2.1 Learning normality-calibrated autoencoder

For n number of input samples with D dimensions $\mathcal{X} = \{x_i\}_{i=1:n}, x \in \mathbb{R}^D$ and the corresponding latent features with d dimensions $\mathcal{Z} = \{z_i\}_{i=1:n}, z \in \mathbb{R}^d$, let an autoencoder is composed of an encoder $f(x) : x \rightarrow z$ and a decoder $g(z) : z \rightarrow \bar{x}$. The general objective of the autoencoder is training f and g to minimise an error between input samples x and the reconstruction results \bar{x} , as follows:

$$\min_{f,g} \mathbb{E}_{x \sim p_{\mathcal{X}}} \|x - \bar{x}\|^2, \quad \bar{x} = g \cdot f(x), \quad (1)$$

where $p_{\mathcal{X}}$ denotes the entire input samples. However, an autoencoder is known to have an over-confidence issue, *i.e.*, low reconstruction error of unseen samples. AD methods using the autoencoder usually identify abnormal samples using the reconstruction error. Therefore, even if the autoencoder takes anomaly samples as inputs, it may not distinguish whether the samples are abnormal or not [Pidhorskyi et al., 2018, Yu et al., 2021]. This over-confidence issue would be more deepened when a training dataset is contaminated.

One straightforward approach to prevent this issue is adding an extra term to maximise reconstruction error for contaminated samples. We define normality-calibrated reconstruction (NCR) loss as follows:

$$\min_{f,g} \mathbb{E}_{x \sim p_{\mathcal{X}^N}} \|x - \bar{x}\|^2 - \mathbb{E}_{x^c \sim p_{\mathcal{X}^C}} \|x^c - \bar{x}^c\|^2, \quad (2)$$

where $p_{\mathcal{X}^N}$ and $p_{\mathcal{X}^C}$ denote the normal samples and contaminated samples, respectively, among input samples $p_{\mathcal{X}}$. Now, we should find out which samples are contaminated to optimise autoencoder using Eq. 2 properly.

2.2 high-confidence normal samples generation using Generative Adversarial Network

We find contaminated samples by using high confident normal samples generated from low entropy latent space. We apply the generative adversarial network (GAN) Goodfellow et al. [2014] framework to do this. The high-confidence normal sample generation via the GAN framework is carried out as follows. Initially, we transform a distribution of all latent features, which are encoded from input samples through the encoder f , to a more knowledgeable probabilistic distribution such as Gaussian distribution. And then, we generate samples using noise signals sampled from the centre of the knowledgeable distribution, *i.e.*, the low entropy space. An adversarial loss for transforming a latent feature distribution to a more knowledgeable probabilistic distribution is defined by follow:

$$\min_f \max_{D_l} \mathbb{E}_{\omega \sim N(\mu_{\mathcal{Z}}, I_d)} [\log D_l(\omega)] + \mathbb{E}_{x \sim P_{\mathcal{X}}} [\log (1 - D_l(f(x)))], \quad (3)$$

where D_l denotes the discriminator for latent features, and $N(\mu_{\mathcal{Z}}, I_d)$ defines a normal distribution with the mean of latent features $\mu_{\mathcal{Z}} \in \mathbb{R}^d$ and a covariance matrix defined by an identity matrix $I_d \in \mathbb{R}^{d \times d}$. $\mu_{\mathcal{Z}}$ is initialised by the mean value of latent features: $\mu_{\mathcal{Z}} = \frac{1}{n} \sum_{i=1}^n z_i$. We would want

67 each component of z to be maximally informative such as each of them to be an independent random
 68 variable. Therefore, the covariance matrix is determined by the $d \times d$ identity matrix.

69 Since as f and D_l are being updated, $\mu_{\mathcal{Z}}$ would be shifted during the training step. $\mu_{\mathcal{Z}}$ is updated at
 70 every training step as follows:

$$\mu_{\mathcal{Z}}^{t+1} = \mu_{\mathcal{Z}}^t - \gamma \frac{1}{m} \sum_{i=1}^m (\mu_{\mathcal{Z}}^t - z_i), \quad \mu_{\mathcal{Z}}^0 = \frac{1}{n} \sum_{i=1}^n z_i^0 \quad (4)$$

71 where $\mu_{\mathcal{Z}}^{t+1}$ and $\mu_{\mathcal{Z}}^t$ denote the $\mu_{\mathcal{Z}}$ on $t + 1$ -th and t -th training step, respectively. m is the batch size
 72 and z_i is i -th latent features on the batch. γ is a learning rate.

73 To generate high confident normal samples, we formulate the following adversarial loss:

$$\min_g \max_{D_s} \mathbb{E}_{x \sim P_{\mathcal{X}}} [\log D_s(x)] + \mathbb{E}_{\omega \sim N(\mu_{\mathcal{Z}}, \sigma I_d)} [\log(1 - D_s(g(\omega))), \quad (5)$$

74 where D_s denotes the discriminator for samples, and $N(\mu_{\mathcal{Z}}, \sigma I_d)$ is a d -dimensional normal distribu-
 75 tion with the mean $\mu_{\mathcal{Z}} \in \mathbb{R}^d$ and the covariance matrix $\sigma I_d \in \mathbb{R}^{d \times d}$. $\mu_{\mathcal{Z}}$ is equivalent to the $\mu_{\mathcal{Z}}$ in
 76 Eq. 3. σI_d is defined by multiplication of a scalar value $\sigma \in [0, 1]$ and the identity matrix I_d . σ is a
 77 hyperparameter to control the compactness of random noise for generating samples using the decoder
 78 g . The smaller σ can give more strong likely to generate high confident normal samples.

79 2.3 Contaminated sample mining and joint learning

80 To predict contaminated samples, we use the generated high confident normal samples as a dictionary.
 81 With the generation process for high confident normal samples: $g(\omega) = \hat{x}$, we construct a latent
 82 feature dictionary $\mathcal{M} = [\hat{z}_i]_{i=1:m}$, $\hat{z}_i = f(\hat{x}_i)$ and $\mathcal{M} \in \mathbb{R}^{m \times d}$, where m is the batch size. By
 83 leveraging \mathcal{M} and given each training batch $\{x_i\}_{i=1:m}$, we define a pseudo contamination score c_i
 84 of each input sample x_i as follow:

$$c_i = \frac{1}{m} \sum_{j=1}^m f(x_i) \cdot \hat{z}_j^T, \quad \hat{z}_j \in \mathcal{M}, \quad (6)$$

85 where T denotes the transpose of the vector. We apply l_2 -normalisation to improve robustness on the
 86 variation of the vector scale of the operation.

87 We predict the contaminated samples by sorting the score in descending order and picking top-
 88 $\tau\%$ samples among the sorted results as the contaminated samples; thus, the number of predicted
 89 contaminated samples are decided by τm that is a multiplication of τ and the batch size m . The
 90 above process is represented as follow:

$$\mathcal{X}^C = \{x_t\}_{t \in C[1:\lceil \tau m \rceil]}, \quad C = \arg \underset{i}{\text{sort}} c_i, \quad w.r.t., 1 \leq i \leq m$$

91 where C is a set of the sorted indices of input batch samples in descending order of the contamination
 92 score (Eq. 6), and \mathcal{X}^C is a set of predicted contaminated samples. $\lceil \cdot \rceil$ denotes the ceiling function.
 93 τ effects of deciding the number of predicted contaminated samples, so it directly affects the AD
 94 performance of our method.

95 The objective function for joint learning the entire components on our method is as follows:

$$\begin{aligned} \min_{f,g} \max_{D_l, D_s} & \underbrace{\mathbb{E}_{x \sim p_{\mathcal{X}^N}} \|x - f \cdot g(x)\|^2 - \mathbb{E}_{x \sim p_{\mathcal{X}^C}} \|x - \bar{x}'\|^2}_{(a)} \\ & + \underbrace{\mathbb{E}_{\omega \sim N(\mu_{\mathcal{Z}}, I_d)} [\log D_l(\omega)] + \mathbb{E}_{x \sim P_{\mathcal{X}}} [\log(1 - D_l(f(x)))]}_{(b)} \\ & + \underbrace{\mathbb{E}_{x \sim P_{\mathcal{X}}} [\log D_x(\omega)] + \mathbb{E}_{\omega' \sim N(\mu_{\mathcal{Z}}, \sigma I_d)} [\log(1 - D_s(g(\omega')))]}_{(c)}, \end{aligned} \quad (7)$$

96 where \bar{x}' is defined by the nearest sample from the given contaminated samples among the generated
 97 high confident normal samples $g(N(\mu_{\mathcal{Z}}, \sigma I_d))$ on the latent feature space. (a), (b), and (c) denote the
 98 NCR loss and the two adversarial losses, respectively.

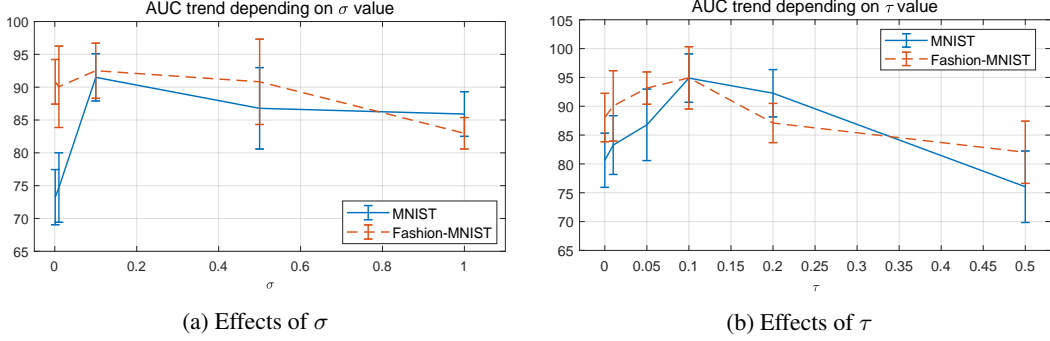


Figure 2: Ablation studies about unsupervised AD performance depending on σ and τ . (a) and (b) represent the trends of AUC with respect to the setting of σ and τ , respectively, on the MNIST and Fashion-MNIST datasets.

99 3 Experiments

100 3.1 Experiment setting and Dataset

101 We follow the unsupervised AD protocol described by Ruff *et al.* [Ruff et al., 2019a]. MNIST and
 102 Fashion-MNIST datasets are used for the experiments. We set one of the classes provided by a dataset
 103 as normal and others as abnormal. After we decide contamination ration $\rho = \frac{A}{N+A}$, where N and A
 104 are the numbers of normal and abnormal samples, respectively, we pick normal samples from the
 105 chosen class and contaminated samples from the remaining classes. In the test phase, the samples of
 106 the normal class are labelled by 0, and other samples are labelled by 1. For the performance analysis,
 107 Receiver Operating Characteristic (ROC) curve and Area Under the Curve (AUC) are used.

108 We employ LeNet-type convolutional neural networks (CNNs) on MNIST and Fashion-MNIST
 109 datasets, where each convolutional module consists of a convolutional layer followed by leaky ReLU
 110 activation functions with leakiness of 0.1. We use the Adam optimiser with the recommended default
 111 hyperparameters [Kingma and Ba, 2015]. The batch size is set to 128. The initial learning rate is 0.01
 112 and decayed every 10 epochs by multiplying 0.1. σ and τ are decided by 0.1 and 0.1 (based on the
 113 results from the ablation study), respectively.

114 3.2 Ablation study

115 We analyse unsupervised AD performance depending on the setting of σ and τ . MNIST and
 116 Fashion-MNIST datasets are used for the ablation study. Ablation studies are conducted based on the
 117 experimental protocol described in the previous section. The contamination ratio ρ is fixed to 0.2.

118 **Parameter analysis on σ :** When σ is too small, then the distribution of sample noise for generating
 119 samples would be too compact so that the generated samples can not provide comprehensive informa-
 120 tion to cover the diverse patterns of normal samples. On the other hand, when σ is too large, then
 121 there is a possibility that the noise can be sampled from low entropy space (*i.e.*, abnormal samples
 122 also can be generated).

123 Figure 2(a) shows the AUC trends depending on the σ . The AUC increases rapidly in the case
 124 of sigma is less than 0.1, and then decreases gradually. This can be interpreted as follows. If the
 125 sampling space is too compact (*i.e.*, when σ is too small), it means that the generated normal sample
 126 does not provide enough information to distinguish the contaminated sample. When sampling space
 127 is too broad (*i.e.*, when σ is too large), it also degrades performance, but the impacts of the broader
 128 sampling space are relatively less than that of the smaller sampling space (e.g., when $\sigma \leq 0.1$). The
 129 best performance is obtained by σ of 0.1.

130 **Parameter analysis on τ :** τ decides the number of predicted contaminated samples per training
 131 batch. The lower τ can provide more precise prediction performance but may not enough to provide
 132 more comprehensive prediction performance. In contrast, when τ is too large, the predicted results
 133 possibly more accurate but also may have a lot of false-positive results.

Table 1: Performance comparison on unsupervised anomaly detection in terms of various contamination ratios ρ . MNIST and Fashion-MNIST datasets are used for the comparison. The **bolded** figures indicate the best performances.

Dataset	ρ	OC-SVM	IF	KDE	CAE	Deep SVDD	SSAD	SS-DGM	Deep SAD	Classification	NCAE
MNIST	.00	96.0±2.9	85.4±8.7	95.0±3.3	92.9±5.7	92.8±4.9	97.9±1.8	92.2±5.6	96.7±2.4	94.5±4.6	94.0±4.2
	.01	94.3±3.9	85.2±8.8	91.2±4.9	91.3±6.1	92.1±5.1	96.6±2.4	92.0±6.0	95.5±3.3	91.5±5.9	97.2±5.2
	.05	91.4±5.2	83.9±9.2	85.5±7.1	87.2±7.1	89.4±5.8	93.4±3.4	91.0±6.9	93.5±4.1	86.7±7.4	97.0±7.1
	.10	88.8±6.0	82.3±9.5	82.1±8.5	83.7±8.4	86.5±6.8	90.7±4.4	89.7±7.5	91.2±4.9	83.6±8.2	92.6±5.7
	.20	84.1±7.6	78.7±10.5	77.4±10.9	78.6±10.3	81.5±8.4	87.4±5.6	87.4±8.6	86.6±6.6	79.7±9.4	89.8±7.4
F-MNIST	.00	92.8±4.7	91.6±5.5	92.0±4.9	90.2±5.8	89.2±6.2	94.0±4.4	71.4±12.7	90.5±6.5	76.8±13.2	91.5±9.7
	.01	91.7±5.0	91.5±5.5	89.4±6.3	87.1±7.3	86.3±6.3	92.2±4.9	71.2±14.3	87.2±7.1	67.3±8.1	94.5 ±4.7
	.05	90.7±5.5	90.9±5.9	85.2±9.1	81.6±9.6	80.6±7.1	88.3±6.2	71.9±14.3	81.5±8.5	59.8±4.6	92.4± 8.2
	.10	89.5±6.1	90.2±6.3	81.8±11.2	77.4±11.1	76.2±7.3	85.6±7.0	72.5±15.5	78.2±9.1	56.7±4.1	91.5±5.7
	.20	86.3±7.7	88.4±7.6	77.4±13.6	72.5±12.6	69.3±6.3	81.9±8.1	70.8±16.0	74.8±9.4	53.9±2.9	88.9±9.2

134 As shown in Figure 2(b), the AUC increases rapidly with τ from 0 to 0.1 and then decreases slowly.
 135 The results can be interpreted as follows. Finding contaminated samples themselves has a large
 136 impact on the AD performance, but the quantity of found samples affects less to AD performance.
 137 But, predicting too many samples may degrade AD performance by taking a great number of false
 138 positives. The best performance is obtained by τ of 0.1.

139 3.3 Comparison with other methods

140 We consider the OC-SVM [Schölkopf et al., 2001], isolation forest (IF) [Liu et al., 2008], and KDE
 141 [Parzen, 1962] for shallow unsupervised baselines. For deep unsupervised competitors, we consider
 142 general binary classifier (supervised), convolutional autoencoders (CAE), deep support vector data
 143 description (Deep SVDD) [Ruff et al., 2018], semi-supervised anomaly detection (SSAD) [Ruff
 144 et al., 2019a], semi-supervised deep generative model (SS-DGM) [Kingma and Ba, 2015], and deep
 145 semi-supervised anomaly detection (Deep SAD) [Ruff et al., 2019a]. We repeat this training set
 146 generation process 10 times per AD setup over all the nine respective anomaly classes and report the
 147 average results over the resulting 90 experiments per contamination ratio.

148 Table 1 shows the quantitative performance comparison depending on the contamination ratio ρ . In
 149 the comparison using the MNIST dataset, the proposed NCAE achieves the best performances except
 150 when the dataset is not contaminated ($\rho = 0.0$). Even compared with semi-supervised approaches
 151 [Ruff et al., 2018, 2019a] which use explicit anomaly samples in the training phase, the NCAE
 152 shows outstanding performances. This trend is also shown in the performance comparison using the
 153 Fashion-MNIST dataset. The NCAE produces the AUC of 91.57 and 88.97 for the Fashion-MNIST
 154 dataset with 0.1 and 0.2 contamination ratios, respectively. Those figures are the best performance
 155 among the listed methods when a dataset is contaminated.

156 The interpretation of the relatively low performance on the uncontaminated dataset ($\rho = 0.0$) is as
 157 follows. Basically, our method is derived under the assumption that a training dataset is contaminated.
 158 Therefore, even if the dataset is not contaminated, the NCAE tries to find some anomaly samples and
 159 maximise the reconstruction errors of the samples during the model training. This process degrades
 160 the performance of our methods as shown in the experimental results. This is a critical defect of our
 161 method;

162 Overall, the comparison results demonstrate the advantage of the proposed NCAE that can detect
 163 anomaly samples on data contamination without prior knowledge or explicit abnormal samples in the
 164 training phase.

165 4 Conclusion

166 In this work, we have proposed NCAE that is a generative method for fully unsupervised anomaly
 167 detection on contaminated data. The experimental results have suggested that the NCAE outperforms
 168 existing methods for fully unsupervised anomaly detection with a large margin, and they have also
 169 provided competitive performances compared with semi-supervised methods using explicit abnormal
 170 samples to train their AD model.

References

- 171
172 S. Akcay, A. Atapour-Abarghouei, and T. P. Breckon. GANomaly: Semi-supervised anomaly
173 detection via adversarial training. In *ACCV*, pages 622–637, 2018.
- 174 A. Berg, J. Ahlberg, and M. Felsberg. Unsupervised learning of anomaly detection from contaminated
175 image data using simultaneous encoder training. *arXiv preprint arXiv:1905.11034*, 2019.
- 176 R. Chalapathy and S. Chawla. Deep learning for anomaly detection: A survey. *arXiv preprint*
177 *arXiv:1901.03407*, 2019.
- 178 J. Chen, S. Sathe, C. C. Aggarwal, and D. S. Turaga. Outlier Detection with Autoencoder Ensembles.
179 In *SDM*, pages 90–98, 2017.
- 180 L. Deecke, R. A. Vandermeulen, L. Ruff, S. Mandt, and M. Kloft. Image anomaly detection with
181 generative adversarial networks. In *ECML-PKDD*, 2018.
- 182 S. M. Erfani, S. Rajasegarar, S. Karunasekera, and C. Leckie. High-dimensional and large-scale
183 anomaly detection using a linear one-class SVM with deep learning. *Pattern Recognition*, 58:
184 121–134, 2016.
- 185 I. Golan and R. El-Yaniv. Deep anomaly detection using geometric transformations. In *NeurIPS*,
186 pages 9758–9769, 2018.
- 187 I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and
188 Y. Bengio. Generative Adversarial Nets. In *NIPS*, pages 2672–2680, 2014.
- 189 N. Görnitz, M. Kloft, K. Rieck, and U. Brefeld. Toward supervised anomaly detection. *Journal of*
190 *Artificial Intelligence Research*, 46:235–262, 2013.
- 191 D. Hendrycks, M. Mazeika, and T. G. Dietterich. Deep anomaly detection with outlier exposure. In
192 *ICLR*, 2019a.
- 193 D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song. Using self-supervised learning can improve
194 model robustness and uncertainty. In *NeurIPS*, pages 15637–15648, 2019b.
- 195 D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015.
- 196 C.-H. Lai, D. Zou, and G. Lerman. Robust subspace recovery layer for unsupervised anomaly
197 detection. In *International Conference on Learning Representations*, 2020. URL [https://](https://openreview.net/forum?id=rylb3eBtwr)
198 openreview.net/forum?id=rylb3eBtwr.
- 199 T. Li, Z. Wang, S. Liu, and W.-Y. Lin. Deep unsupervised anomaly detection. In *Proceedings of the*
200 *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3636–3645, 2021.
- 201 F. T. Liu, K. M. Ting, and Z.-H. Zhou. Isolation Forest. In *ICDM*, pages 413–422, 2008.
- 202 Y. Liu and Y. F. Zheng. Minimum enclosing and maximum excluding machine for pattern description
203 and discrimination. In *ICPR*, pages 129–132, 2006.
- 204 G. Pang, C. Shen, and A. van den Hengel. Deep anomaly detection with deviation networks. In *KDD*,
205 pages 353–362, 2019.
- 206 E. Parzen. On Estimation of a Probability Density Function and Mode. *The annals of mathematical*
207 *statistics*, 33(3):1065–1076, 1962.
- 208 S. Pidhorskyi, R. Almhosen, and G. Doretto. Generative probabilistic novelty detection with adver-
209 sarial autoencoders. In *NeurIPS*, pages 6822–6833, 2018.
- 210 L. Ruff, R. A. Vandermeulen, N. Görnitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and
211 M. Kloft. Deep one-class classification. In *ICML*, volume 80, pages 4390–4399, 2018.
- 212 L. Ruff, R. A. Vandermeulen, N. Görnitz, A. Binder, E. Müller, K.-R. Müller, and M. Kloft. Deep
213 semi-supervised anomaly detection. *arXiv preprint arXiv:1906.02694*, 2019a.

- 214 L. Ruff, Y. Zemlyanskiy, R. Vandermeulen, T. Schnake, and M. Kloft. Self-attentive, multi-context
215 one-class classification for unsupervised anomaly detection on text. In *ACL*, pages 4061–4071,
216 2019b.
- 217 B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the Support
218 of a High-Dimensional Distribution. *Neural computation*, 13(7):1443–1471, 2001.
- 219 H. Song, Z. Jiang, A. Men, and B. Yang. A hybrid semi-supervised anomaly detection model for
220 high-dimensional data. *Computational Intelligence and Neuroscience*, 2017.
- 221 J. Wang, P. Neskovic, and L. N. Cooper. Pattern classification via single spheres. In *International
222 Conference on Discovery Science*, pages 241–252. Springer, 2005.
- 223 J. Yu, Y. Lee, K. C. Yow, M. Jeon, and W. Pedrycz. Abnormal event detection and localization via
224 adversarial event prediction. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- 225 S. Zhai, Y. Cheng, W. Lu, and Z. Zhang. Deep structured energy based models for anomaly detection.
226 In *ICML*, volume 48, pages 1100–1109, 2016.
- 227 B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen. Deep autoencoding
228 gaussian mixture model for unsupervised anomaly detection. In *International Conference on
229 Learning Representations*, 2018.