

Feature relevance XAI in anomaly detection: reviewing approaches and challenges

Julian Tritscher^{1,*}, Anna Krause¹ and Andreas Hotho¹

¹ Data Science Chair, University of Würzburg, Würzburg, Germany

Correspondence*:

Data Science Chair, University of Würzburg, Sanderring 2, 97074 Würzburg, Germany
tritscher@informatik.uni-wuerzburg.de

2 ABSTRACT

3 With complexity of artificial intelligence systems increasing continuously in past years, studies
4 to explain these complex systems have grown in popularity. While much work has focused on
5 explaining artificial intelligence systems in popular domains such as classification and regression,
6 explanations in the area of anomaly detection have only recently received increasing attention
7 from researchers. In particular, explaining singular model decisions of a complex anomaly
8 detector by highlighting which inputs were responsible for a decision, commonly referred to as
9 local post-hoc feature relevance, has lately been studied by several authors. In this paper, we
10 systematically structure these works based on their access to training data and the anomaly
11 detection model, and provide a detailed overview of their operation in the anomaly detection
12 domain. We demonstrate their performance and highlight their limitations in multiple experimental
13 showcases, discussing current challenges and opportunities for future work in feature relevance
14 XAI for anomaly detection.

15 **Keywords:** explainable artificial intelligence, xai, feature relevance, anomaly detection, artificial intelligence, machine learning, review

1 INTRODUCTION

16 Within the last years, artificial intelligence (AI) systems have transformed from simple and interpretable
17 decision systems to complex and highly opaque architectures that are commonly comprised of millions of
18 parameters (Arrieta et al., 2020). With increasing deployment of these highly performing opaque AIs in
19 practice, many application areas have identified a need for explaining the reasoning of complex AI systems.
20 Motivations for explaining these systems range from reducing manual inspection efforts in domains such
21 as medicine (Tjoa and Guan, 2021), to legal requirements for AIs that significantly affect users (Goodman
22 and Flaxman, 2017). As a result, explainable AI (XAI) has become a popular area of research. While the
23 field itself has a longer history with several early applications (Setiono and Leow, 2000; Féraud and Clérot,
24 2002; Robnik-Šikonja and Kononenko, 2008), a lot of research has been conducted in the last six years to
25 provide explanations mainly for common AI tasks such as classification and regression problems (Arrieta
26 et al., 2020). In the area of anomaly detection, research on explainability has taken off more recently,
27 motivated through use in critical security applications such as intrusion and fraud detection (Antwarg
28 et al., 2021), and the desire to decrease manual investigation efforts by domain experts that inspect found
29 anomalies (Sipple and Youssef, 2022).

30 With the increasing interest on explaining anomaly detection within recent years, first works have started
31 to categorize this emerging research field. While Sejr and Schneider-Kamp (2021) discuss the process of
32 explaining anomaly detection from a user perspective, Nonnenmacher et al. (2022) aggregate anomaly
33 detection XAI work that was specifically designed for tabular data. Panjei et al. (2022) and Yepmo et al.
34 (2022) both provide a general overview of the field of anomaly XAI that categorizes the general types
35 of explanations that may be used to explain anomaly detectors, splitting XAIs by the granularity of their
36 given outputs. Panjei et al. (2022) discuss explanations that return a ranking of found anomalies, XAIs
37 that find causal interactions of outliers, and methods that find relevant features. They focus largely on
38 white box models that find characteristics of outliers in big data. Yepmo et al. (2022) provide an illustrated
39 introduction to four general types of anomaly explanations, e.g. ones that return relevant features or
40 decision rules, and name representative approaches. The authors discuss limitations of the general types
41 of explanations only at a high level, without distinguishing between different approaches. In contrast, we
42 focus on reviewing one specific type of anomaly explanation in-depth. This focused view allows us to
43 construct a fine-grained systematic categorization of different algorithmic approaches and investigate each
44 algorithm in detail. Our review highlights low level limitations of XAI algorithms in anomaly detection
45 that constitute relevant areas for future work.

46 In this work, we provide an in-depth review of approaches that produce explanations commonly referred
47 to as *local post-hoc feature relevance XAIs* (Arrieta et al., 2020) in the field of anomaly detection. While a
48 variety of XAIs exist that yield different types of explanations as output, *feature relevance XAIs* explain the
49 decision process of anomaly detection models through highlighting relevant input features, providing as
50 output a relevance score for each input feature. They constitute the currently most used type of explanation
51 in anomaly detection (Yepmo et al., 2022). Applying feature relevance XAIs in a *local* fashion, i.e. per data
52 point, results in highlighting relevant input features that lead an anomaly detection model to identifying a
53 singular data point as an anomaly, in contrast to XAIs that provide a global explanation of general model
54 behavior. This provides additional information regarding a singular found anomaly to manual investigators
55 and reduces their inspection efforts. Further, in contrast to ante-hoc approaches that describe inherently
56 explainable anomaly detectors such as simple linear models, *post-hoc XAIs* describe dedicated XAI
57 approaches that are applied to already fully trained anomaly detectors, allowing the use of highly complex
58 and well performing model architectures without constraining their complexity during model training. The
59 resulting sub-field of local post-hoc feature relevance XAI, which we will refer to in abbreviated form as
60 *feature relevance XAI* in the remaining paper, has recently received increasing attention within the domain
61 of anomaly detection. We systematically review approaches from this sub-field that have been applied to
62 anomaly detection in the remaining paper.

63 We provide a structured characterization of the reviewed approaches in Figure 1, where we group
64 approaches based on their reliance on training data on one hand, and on the anomaly detection architecture
65 on the other hand. This categorization leads from completely model-agnostic XAI approaches that utilize
66 additional assumptions and information obtained through training data, to model-specific XAIs that heavily
67 rely on the model structure to obtain feature relevance explanations. Additionally, we identify two groups
68 of hybrid approaches that access both the model and information regarding the underlying data. While
69 perturbation-based approaches restrict their model access purely to allowing inference on data that may
70 be augmented according to data assumptions, gradient-based approaches require access to the first order
71 derivatives of differentiable anomaly detection models. We provide an in-depth introduction to these groups

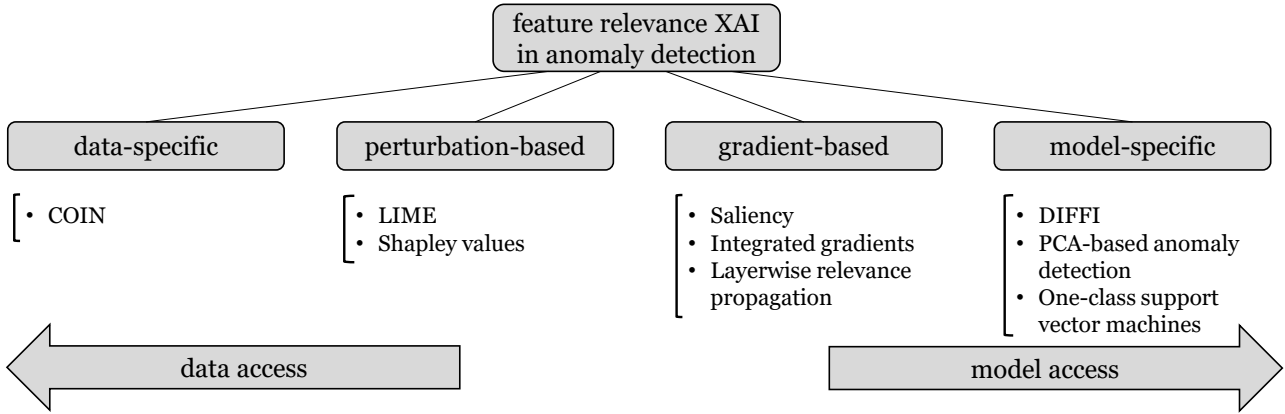


Figure 1. Overview of the reviewed feature relevance approaches in anomaly detection, structured by their use of information from data and from the underlying anomaly detection model.

72 of approaches, and demonstrate their limitations and challenges in multiple showcases to outline future
 73 research directions of feature relevance XAI in anomaly detection.¹

74 The remainder of the paper is structured as follows: Section 2 formally introduces the tasks of anomaly
 75 detection and feature relevance explanations, as well as the data, model architectures, and performance
 76 metrics we use in our showcases. Section 3 covers data-specific approaches that possess no access to
 77 the anomaly detection model, instead generating their explanations through training data. Section 4
 78 introduces perturbation-based approaches that generate explanations through repeatedly querying the
 79 anomaly detection model with altered, so called *perturbed* data points. Section 5 includes gradient-based
 80 approaches that require differentiability of the anomaly detection model and utilize gradients that contain
 81 knowledge of the inner model structure to obtain explanations. Section 6 presents model-specific approaches
 82 that are developed for specific model architectures and take full advantage of the model structure to generate
 83 their explanations. In Section 7 we conclude by discussing the overarching limitations of feature relevance
 84 explanations in anomaly detection and highlighting future research areas within the domain.

2 METHODOLOGY

85 Before we review existing feature relevance XAI approaches, we briefly define the tasks of anomaly
 86 detection and feature relevance XAI, as well as give a brief overview of the data, anomaly detection models,
 87 and XAI evaluation metrics we use to showcase XAI approaches and their limitations throughout this study.

88 2.1 Anomaly Detection

89 Anomaly detection, as laid out by Chandola et al. (2009), describes the task of identifying anomalous
 90 behavior in data that contains well-defined normal behavior. For data points $x \in \mathcal{X} \subseteq \mathbb{R}^d$ of dimensionality
 91 d , anomaly detection is the identification of anomalous data points through a model $m(x)$, where $m(x)$
 92 may be modeled as a binary classification, a probabilistic estimation of anomalies, or as a regression task
 93 that assigns each point an anomaly score. In this work, we view anomaly detection as regression task
 94 $m(x) : \mathcal{X} \rightarrow [0, \text{inf}]$ with lower scores representing normal data and higher scores for anomalies.

95 While anomaly detection at a high level is only a subset of classification or regression, the unique
 96 challenges in anomaly detection arise from specific data characteristics: only the normal behavior in

¹ Code for our showcases is available at <https://professor-x.de/feature-relevance-AD>

97 anomaly detection is well-defined and normal data is typically readily available, but anomalies may vary
98 greatly in behavior with only a small number of anomalies that are known during training. As a result,
99 proposed approaches typically focus on the well-defined normal data to be able to identify potentially
100 unseen types of anomalous behavior, e.g. through encircling observed normal behavior in one-class
101 support vector machines (Schölkopf et al., 2001), assessing the density around data points in kernel density
102 estimation (Terrell and Scott, 1992), or learning a reconstruction of the normal behavior with autoencoder
103 neural networks (Goodfellow et al., 2016).

104 2.2 Feature Relevance Explanations

105 Local post-hoc feature relevance explanations explain the model prediction $m(x)$ for a specific input x
106 through assigning a score to each input feature, creating an explanation $f(x, m) \in \mathbb{R}^d$ that reflects how
107 much each input feature influenced the final prediction according to model m . In the domain of anomaly
108 detection, feature relevance explanations are commonly applied to anomalous data points, and focus on
109 highlighting the relevant features that lead the anomaly detection model to identify the data point as an
110 anomaly (Yepmo et al., 2022).

111 2.3 Data

112 While there is no shortage of datasets for anomaly detection, most of these do not include ground truth
113 explanations for anomalies. Since this ground truth enables an otherwise challenging direct comparison
114 and quantitative judgment of explanations generated by XAI approaches, we select two datasets for the
115 showcases conducted in this review that offer these ground truth explanations: MVTec (Bergmann et al.,
116 2019) and ERP (Tritscher et al., 2022a).

117 MVTec (Bergmann et al., 2019) is an anomaly detection dataset for industrial visual fault detection. The
118 dataset contains 15 texture and object classes with the training set for each category containing only normal
119 images, e.g., without defects, and the test set containing images with defects and without defects. The
120 defects are annotated with manually created ground truth pixel maps, with binary indications of pixels that
121 are part of the defect. The dataset has been previously used to evaluate feature relevance XAI approaches
122 by Ravi et al. (2021), although their evaluations are limited to qualitative inspections of results. To instead
123 generate quantitative results of XAI performance, we use the ground truth anomaly segmentation maps as
124 ground truths for explanations. For our showcases, we focus on the *grid* class from the dataset that contains
125 264 high resolution images of normal wire mesh for training and 57 images with different faults and ground
126 truth for testing. We choose this class as it has the highest detection accuracy of the used anomaly detection
127 model. This limits the influence of poor model performance on the quality of the obtained explanations,
128 which we motivate further in Section 2.4.

129 ERP (Tritscher et al., 2022a) is a synthetic enterprise resource planning (ERP) dataset generated by using
130 a serious game within a real ERP system (Léger et al., 2007). The data includes financial documents from
131 a simulated production company, where different financial fraud scenarios have been committed within
132 the simulation. Additionally, the provided fraudulent data points come with ground truth features that are
133 indicative of the fraud case according to auditing experts, which we utilize as ground truth explanations.
134 For analysis in this work, we rely on the joint machine-learning ready data provided by Tritscher et al.
135 (2022a) that focuses on the financial accounting data. We utilize their run *normal 2* that contains 32337
136 data points of purely normal operation for training the anomaly detector and evaluating explanations on the
137 86 different fraud cases contained in their run *fraud 3*. We choose these runs following the experimental

138 setup of Tritscher et al. (2022b), again using *fraud 3* as the dataset with highest performance of the used
139 anomaly detection model and the corresponding normal behavior of *normal 2*.

140 2.4 Models

141 To showcase XAI algorithms on the introduced data, we select an anomaly detection model with high
142 detection performance through common metrics such as AUC-PR and AUC-ROC scores from literature
143 for each dataset. We specifically require high performance from our anomaly detectors to not obscure the
144 quantitative XAI evaluation. With poorly performing models a miss match of ground truth and explanation
145 may be caused by the model, and not just the XAI approach, preventing the result from reflecting the XAI
146 performance.

147 For the MVTEC image dataset, Kauffmann et al. (2020b) train kernel density estimation (Rosenblatt, 1956),
148 deep support vector data description (Ruff et al., 2018), and autoencoder neural networks (Goodfellow et al.,
149 2016) on MVTEC data. While their models show high anomaly detection performance, further analyzes by
150 the authors reveal that their models and model ensembles use spurious correlations in the data, which may
151 skew a quantitative XAI evaluation. Wang et al. (2021) propose a student-teacher neural network that is
152 designed for segmenting anomalous regions within the MVTEC images. The network incorporates a teacher
153 network that consists of three pre-trained feature extraction layers from the popular ResNet-18 architecture
154 (He et al., 2016) and a randomly initialized student that possesses the same network architecture as the
155 teacher and is trained to mimick the pre-trained teacher on normal training data. While the resulting
156 student-teacher architecture directly outputs image segmentation maps with highlighted anomalous regions,
157 it can be adapted to image-level anomaly detection through adding a mean pooling step to the final output.
158 This creates a well-performing image-level anomaly-detector that is capable of finding anomalies within
159 the MVTEC data both on an image- and a pixel-level and can be used as a test-bed for the investigated XAIs.

160 For the ERP dataset, Tritscher et al. (2022b) conduct a hyperparameter study of multiple anomaly
161 detectors on the data, finding architectures that yield good results on the dataset. For our showcases, we
162 select their second best performing model, the autoencoder neural network (Goodfellow et al., 2016)
163 architecture, with their found hyperparameters as they show that their best performing one-class support
164 vector machine (Schölkopf et al., 2001) exhibits an erratic decision process that may influence a quantitative
165 XAI evaluation and autoencoder networks are commonly studied in the domain of explainable anomaly
166 detection Ravi et al. (2021); Antwarg et al. (2021); Müller et al. (2022).

167 2.5 Evaluation Metrics

168 To showcase the performance of different feature relevance XAI approaches, we utilize the binary ground
169 truth explanations contained in the datasets that denote for each input feature whether the feature was
170 indicative of the underlying anomaly (1) or part of normal behavior (0). To generate quantitative results
171 with this type of ground truth explanation, a performance metric for comparing ground truth with generated
172 explanations is required.

173 Hägele et al. (2020) use the well known area under the receiver operating characteristic (ROC) as metric
174 for their feature relevance evaluation on medical image data. As ROC scores are calculated using the true
175 positive rate over increasing threshold values, early true positives are more impactful to the resulting area
176 under the curve. When applied to feature relevance, this corresponds to a stronger focus on finding truly
177 relevant features within the top scoring features of a given explanation. This is an intuitive metric, as
178 anomaly detectors do not need to identify all anomalous features within an anomaly, but may sufficiently
179 detect the anomaly by focusing heavily on few features that are indicative of the anomalous behavior.

180 To complement the ROC score, we also report cosine similarity (COS) as used for feature relevance
 181 evaluation by Kauffmann et al. (2020b), which reflects the similarity of the found feature relevance
 182 explanations to the entire ground truth. Intuitively, this corresponds to how well an obtained explanation
 183 finds all truly anomalous features. This metric also holds interesting properties in the the case of non-binary
 184 ground truths, since COS respects the magnitudes of the ground truth feature relevance.

185 Both metrics can be calculated for each data point individually, and can then be aggregated across
 186 multiple anomalous data points. In this work, we therefore report mean and standard deviation of the
 187 resulting metrics across all anomalies.

3 DATA-SPECIFIC EXPLANATIONS

188 Data-specific explanations identify relevant feature values of anomalies entirely through training data
 189 without any access to the anomaly detection model. The anomalies themselves are found by an anomaly
 190 detection model, effectively making data-related explanations post-hoc XAIs. However, these approaches
 191 act independently of the anomaly detection model and identify relevant features in given anomalies entirely
 192 through their own assumptions.

193 3.1 Contextual Outlier Interpretation

194 Contextual Outlier INterpretation (COIN) (Liu et al., 2018), to our knowledge currently the only data-
 195 specific post-hoc feature relevance XAI approach, explains an anomalous data point x found by an anomaly
 196 detection model m by determining how much it’s input features are responsible for separating x from
 197 training data \mathcal{X}_{train} . As a first step, COIN extracts context data points \mathcal{C} from the normal data within
 198 \mathcal{X}_{train} that are close to x in feature space through nearest neighbors such that $nn(x, \mathcal{X}_{train}) = \mathcal{C}$. Since
 199 several distinct types of normal behavior might exist in the data, COIN then uses clustering $cl(\mathcal{C}, c) = \mathcal{C}_c$ to
 200 separate the context data points \mathcal{C} into individual groups c with similar behavior. For each of these groups,
 201 a decision boundary separating \mathcal{C}_c from the anomaly x is learned via a linear support vector machine s
 202 (Boser et al., 1992) with loss $\mathcal{L}_s(x, \mathcal{C}_c)$ and an L1 regularization term $\Omega(s)$ through

$$S_c(x) = \arg \min_s \mathcal{L}_s(x, \mathcal{C}_c) + \Omega(s). \quad (1)$$

203 Letting $w_c \in \mathbb{R}^d$ denote the weights of the resulting linear support vector machine S_c for context group
 204 c , the relevance of individual feature values within x are then obtained through the weights of the SVM
 205 through

$$f_c(x_i) = abs(w_{c,i})/\gamma_{c,i}, \quad (2)$$

206 where $\gamma_{c,i}$ denotes the average distance between data points in \mathcal{C}_c for the i th feature. To obtain the final
 207 feature relevance scores of anomaly x , the feature relevance scores of individual context groups are
 208 averaged. This results in the following process for feature relevance explanations:

$$f_{COIN}(x_i, \mathcal{X}_{train}) = (1/|nn(x, \mathcal{X}_{train})|) \sum_c |cl(nn(x, \mathcal{X}_{train}), c)| \cdot f_c(x_i). \quad (3)$$

209 3.2 Limitations

210 Data-specific feature relevance XAIs explain found anomalies purely from the data domain, and are
 211 therefore applicable without any access to the anomaly detection model. Due to this complete separation

212 of anomaly detection model and explanation approach, the XAI needs to build its feature relevance
 213 explanations purely relying on given data. As observed in the introduced COIN framework, this requires
 214 additional assumptions regarding the data in multiple steps during the explanation process. Since COIN
 215 relies both on a nearest neighbors algorithm to identify the local context data points around a given
 216 anomaly, and on clustering to separate multiple distinct types of normal behavior in the data, this requires
 217 the definition of a meaningful distance function within the data. Obtaining reasonable assumptions regarding
 218 the distance metric for a given dataset is a non-trivial task, effectively requiring the construction of an
 219 additional, well-performing, distance-based anomaly detection system to obtain high quality explanations
 220 for a given dataset. As a result, if such a well-performing distance-based anomaly detection system is
 221 not available, e.g. in domains where distance-based anomaly detectors perform poorly in general, an
 222 application of the COIN framework may yield poor results due to it’s internal reliance on the construction
 223 of an additional anomaly detector.

4 PERTURBATION-BASED EXPLANATIONS

224 In contrast to data-driven approaches that only access the final decision of an anomaly detection model
 225 $m(x)$ for a given anomalous data point x , perturbation approaches allow free access of the model decision
 226 function m on arbitrary data points. While this does not provide direct knowledge on the structure of the
 227 anomaly detection model, effectively treating m as a black box, it provides an opportunity to probe the
 228 model behavior. Perturbation approaches use the access to the anomaly detection function m by repeatedly
 229 constructing synthetic data points x' through altering the given anomalous data point x , and probing the
 230 anomaly detection model’s reaction to the alterations by applying the model to the synthetic data points
 231 through $m(x')$.

232 To obtain relevance scores for individual features, this probing procedure is used to remove features and
 233 feature combinations from the anomaly x and measure the model’s reaction to the presence and absence
 234 of features. Perturbation approaches alter an anomalous data point $x = [x_1, x_2, \dots, x_d]$ of dimensionality
 235 d by determining a set of features $K \subseteq \{1, 2, \dots, d\}$ to keep, and subsequently deleting, i.e. perturbing,
 236 the K^C remaining features not in K from data point x , where K^C denotes the complement of K (i.e.
 237 $K^C = \{1, 2, \dots, d\} \setminus K$). This perturbation procedure is used by several XAIs repeatedly on a single data
 238 point x to gather information on the behavior of the machine learning model when specific feature values
 239 within x are removed, allowing them to identify single features and feature groups that determine the model
 240 output. Since a large amount of machine learning models are not capable of handling missing values, the
 241 construction of perturbed data points is commonly achieved not through deletion but through replacing the
 242 values in K^C with additional reference data $r \in \mathbb{R}^d$ through $h(x, r, K) = [x_K, r_{K^C}]$.

4.1 Local Interpretable Model-Agnostic Explanations

244 Local Interpretable Model-agnostic Explanations (LIME) (Ribeiro et al., 2016) generates explanations
 245 for model decisions on single data points x through the perturbation procedure. LIME generates a synthetic
 246 dataset \mathcal{X}' around anomaly x through $s : \mathbb{R}^d, \mathbb{R}^d \rightarrow \mathcal{X}'$ by perturbing x with reference data r and sampling
 247 the features to perturb from a uniform distribution such that

$$s(x, r) = \mathcal{X}' \sim U(\{h(x, r, K), K \subseteq \{1, 2, \dots, d\}\}). \quad (4)$$

248 These synthetic data points are then weighted through a proximity measure π_x that indicates the proximity
 249 of the synthetic points to the original data point x to explain. Using this synthetic data, an explanation is

250 then obtained through the parameters of a linear and therefore interpretable model with linear coefficients
 251 $w \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$, that is trained to mimick the original model m on the synthetic data points \mathcal{X}'
 252 in the proximity π_x through a loss function $\mathcal{L}(m, w, b, \mathcal{X}', \pi_x)$. This linear model is regularized through
 253 a complexity measure $\Omega(w)$, which enforces simple and readily interpretable linear coefficients w . As a
 254 result, LIME generates explanations for a data point x by linearly approximating the original model m in
 255 the local proximity π_x through

$$(W(x, m, r), B(x, m, r)) = \arg \min_{(w, b)} \mathcal{L}(m, w, b, s(x, r), \pi_x) + \Omega(w). \quad (5)$$

256 This results in a local linear model with one linear coefficient for each input feature. As a result, the linear
 257 coefficients show the relevance of each feature in the local vicinity of \mathcal{X}' and can be taken directly as
 258 feature relevance explanations through

$$f_{LIME}(x, m, r) = W(x, m, r). \quad (6)$$

259 Ravi et al. (2021) directly apply LIME on the anomaly detection MVTEC dataset with a brief qualitative
 260 demonstration of results. Further, Zhang et al. (2019) apply LIME on multiple anomaly detection datasets
 261 from the security domain that focus on intrusion and malware detection. While LIME yields both positive
 262 and negative contributions to the model output, Zhang et al. (2019) only retain contribution signals that
 263 cause an increased anomaly score. They also introduce an additional, optional loss term based on KL
 264 divergence that allows for determining the desired distribution of output explanation scores.

265 4.2 Shapley Value Explanations

266 The Shapley value (Shapley, 1997), a well-known result from cooperative game theory, describes a unique
 267 solution to fairly distributing cooperatively achieved gain among n cooperating players by measuring the
 268 achieved gain of partial coalitions. The solution provided by Shapley uniquely satisfies desirable fairness
 269 properties such as permutation in-variance of coalitions and zero gain for players not included in the
 270 coalition, among others. The Shapley value ϕ_i for a single player i represents the gain generated by player i
 271 and can be computed through iteratively measuring the gain of all coalitions without player i in comparison
 272 to the same coalition with player i included, giving

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (v(S \cup \{i\}) - v(S)) \quad (7)$$

273 for the set of all players $N = \{1, 2, \dots, n\}$ and a function $v(S)$ to compute the gain of a coalition S .

274 Applying Shapley values to the domain of feature relevance explanations, as done by Lundberg and Lee
 275 (2017), is achieved by viewing the features of x as players, building coalitions through perturbations, i.e.
 276 through keeping and replacing features, and computing the gain as the outcome of applying the model on
 277 the synthetic data point from the coalition, giving

$$f_{Shapley}(x_i, m, r) = \sum_{K \subseteq N \setminus \{i\}} \frac{|K|!(d - |K| - 1)!}{d!} (m(h(x, r, K \cup \{i\})) - m(h(x, r, K))). \quad (8)$$

Table 1. Mean and standard deviation of perturbation XAI performance comparing to ground truth explanations over all anomalies for ERP and MVTec data respectively.

(A) ERP			(B) MVTec		
XAI	ROC	COS	XAI	ROC	COS
noise	22.7 (7.0)	-28.7 (5.3)	noise	50.2 (1.8)	6.3 (2.7)
noise×input	52.3 (13.0)	-27.2 (8.2)	noise×input	32.3 (10.7)	-4.5 (5.5)
LIME	75.7 (3.9)	28.6 (8.3)	LIME	56.6 (9.2)	10.9 (13.7)
SHAP	74.4 (17.1)	32.3 (25.1)	SHAP	64.5 (19.9)	-5.3 (2.8)

278 Since computing the true Shapley value as feature relevance is prohibitively resource-intensive for
 279 reasonably sized numbers of features d , multiple approaches exist for estimating Shapley values. As the
 280 predominant work in XAI, SHapley Additive exPlanations (SHAP) (Lundberg and Lee, 2017) shows that
 281 proposing slight alterations to existing XAI approaches can yield approximate Shapley value explanations.
 282 For their approach "kernel-SHAP", the authors adapt the perturbation framework of LIME, showing that
 283 LIME is capable of recovering an approximation of Shapley values using the following choices of proximity
 284 kernel π_x and regularization term $\Omega(g)$ when fitting LIME's linear approximation model in Equation (5):

$$\pi_x = \frac{d-1}{\binom{d}{K} \cdot K \cdot (d-K)}, \quad \Omega(g) = 0 \quad (9)$$

285 For datasets with high dimensionality d and a known hierarchy between dimensions (e.g. local dependencies
 286 in images), "partition-SHAP" extends this approach to groups of features through the game-theoretic
 287 extensions to Owen values (Owen, 1977) and achieves faster run times as a result.

288 Shapley value explanations are some of the most used approaches in anomaly detection, with multiple
 289 applications on reconstruction-based anomaly detectors such as autoencoder neural networks (Ravi et al.,
 290 2021; Antwarg et al., 2021; Tritscher et al., 2022b; Müller et al., 2022). While Ravi et al. (2021); Tritscher
 291 et al. (2022b) apply Shapley value estimation directly on the final anomaly score of the reconstruction-based
 292 anomaly detection model, Antwarg et al. (2021) first identify the features with highest reconstruction errors
 293 and apply kernel-SHAP directly on the most deviating features. Müller et al. (2022) further extend this
 294 approach to categorical one-hot encoded data by averaging over groups of one-hot encoded features.

295 4.3 Showcase and Limitations

296 4.3.1 Showcase of perturbation approaches

297 To be able to discuss the application of perturbation approaches to anomaly detection and showcase
 298 the resulting limitations in detail, we first demonstrate the performance of the two previously introduced
 299 approaches LIME and SHAP using their default parameter settings on the datasets MVTec and ERP
 300 described in Section 2.3. While we use kernel-SHAP for all applications of SHAP on ERP data, we use the
 301 authors' partition-SHAP implementation for the large image dataset of MVTec to maintain computational
 302 feasibility. We compute feature relevance explanations on all anomalies in the respective test datasets and
 303 compare the resulting explanations to ground truth using the ROC score and cosine similarity as discussed
 304 in Section 2.5. Additionally, to ease interpretability of results, we introduce two random baselines that
 305 include explanations sampled from random uniform noise, as well as a multiplication of random uniform
 306 noise with the anomalous input itself (noise×input). Table 1A shows that both LIME and SHAP are capable
 307 of highlighting relevant features on ERP data, demonstrating considerably higher scores than random noise.

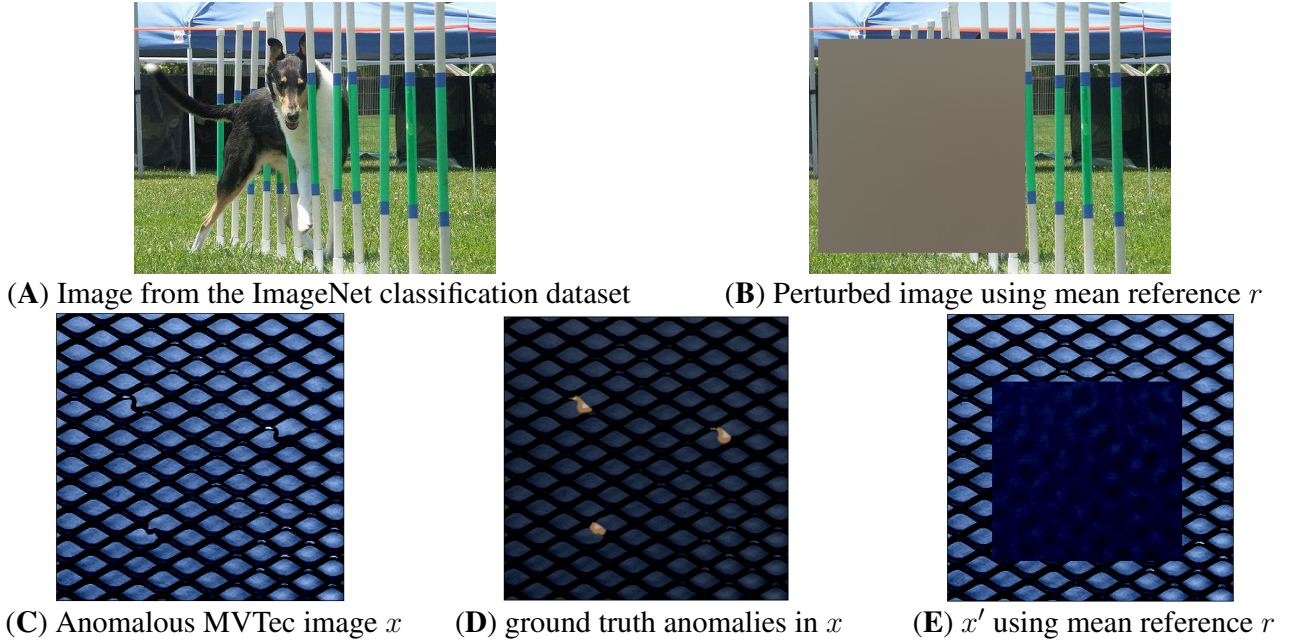


Figure 2. Demonstration of perturbation with mean reference r in classification (A-B) and anomaly detection (C-E). In classification, mean reference is capable of completely removing the class signal 'dog' in image B. In anomaly detection, replacing all areas that contain anomalies (highlighted in image D) with mean reference r introduces new anomalous signals in the resulting data point in image E.

308 On the image data of MVTec in Table 1B, however, both approaches perform poorly on all metrics with
 309 only small improvements over the random baselines.

310 4.3.2 Limitation: choice of reference values r

311 One key aspect of perturbation-based explanation approaches is the choice of reference data r for
 312 removing signal and representing missing information, which is a non-trivial question that is still unsolved
 313 in current research (Ancona et al., 2019). Common references that stem from well researched tasks such as
 314 image classification include replacing feature values with zero values or averages obtained from training
 315 data (used by LIME and SHAP as default in Table 1). We demonstrate this on an image of a dog in
 316 Figure 2A, taken from the well known ImageNet classification dataset (Russakovsky et al., 2015). When
 317 classifying a dog within the image, perturbing features through mean values from data as reference r (here
 318 calculated from the ImageNet validation split as demonstration), intuitively removes any signal present
 319 in the replaced features that might be indicative of the dog (see Figure 2B). As a result, mean values are
 320 capable of removing the relevant signal on perturbed input features in this setting.

321 Within the domain of anomaly detection, however, these fixed reference values might introduce unwanted
 322 signals into the data. We demonstrate this on an anomalous data point from the MVTec test dataset that
 323 contains a bent wire anomaly in an otherwise normal wire mesh (Figure 2C). The anomalous inputs
 324 according to the ground truth explanations are highlighted in Figure 2D. Replacing a region that covers
 325 all anomalous inputs with mean values from the MVTec training data may still yield an anomalous data
 326 point x' that does not represent the well-defined normal behavior of a wire mesh (see Figure 2E). Even
 327 though all inputs that contain anomalous entries have been replaced from the initial anomaly, the resulting
 328 image may still be declared as anomaly by the model and therefore prevent XAI approaches from finding
 329 the relevant anomalous inputs.

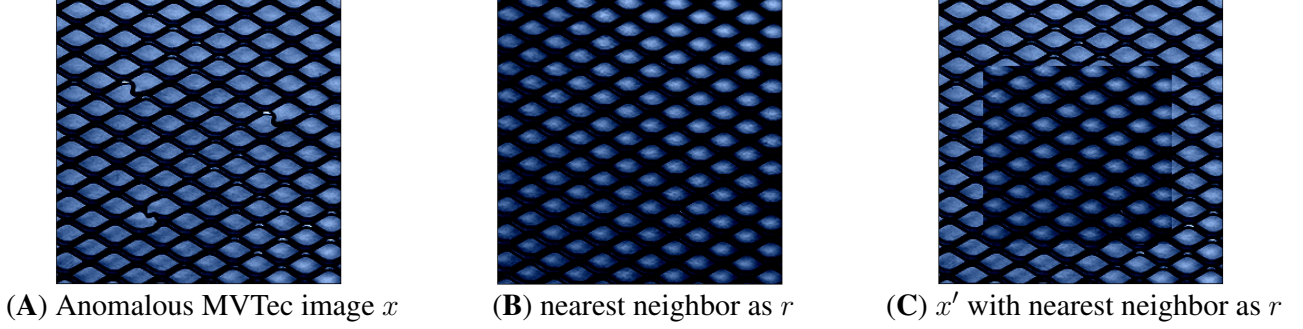


Figure 3. Demonstration of perturbation with nearest neighbor in normal train data as reference r : while r is visually closer to the anomaly x than the mean of training data from Figure 2E, the perturbed point x' still shows highly anomalous characteristics on the replacement borders.

330 4.3.3 Finding optimal reference values r in anomaly detection

331 To alleviate this issue, reference values r have in the past been chosen in the context of the data point x ,
 332 e.g. through finding nearest neighbors to x within normal training data that is both similar to x and lies
 333 within the normal data manifold (Takeishi and Kawahara, 2020). While this can indeed produce normal
 334 data points after perturbation for some groups of retained feature values K (e.g. when replacing all values
 335 within x), for some values of K the combination of anomalous data point x and it’s nearest neighbor might
 336 still introduce further unwanted anomalies as visualized in Figure 3.

337 To achieve better perturbation-based explanations, Takeishi and Kawahara (2020) propose to find r
 338 dynamically dependent on the data point x and features to keep K . To additionally ensure that the perturbed
 339 features make the resulting data point more normal, Takeishi and Kawahara (2020) generate the synthetic
 340 data point x'_{opt} by minimizing the model output in the local neighborhood \mathcal{N}_x of the original data point
 341 while constraining the features in K to their original values in x , giving

$$x'_{opt} = \arg \min_{\hat{x} \in \mathcal{N}_x} m(\hat{x}) \quad s.t. \quad \hat{x}_i = x_i, \forall i \in K. \quad (10)$$

342 Takeishi and Kawahara (2020) further relax this generation procedure by searching for a local minimum of
 343 Equation (10) instead through

$$x'_{lopt} = \arg \min m(\hat{x}) + \gamma \cdot dist(x, \hat{x}) \quad (11)$$

344 using a distance function $dist : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, which may be minimized through constrained optimization
 345 with the constraints $\hat{x}_i = x_i, \forall i \in K$. To further reduce the computational overhead required for synthetic
 346 data generation on data points with reasonably low dimensionality ($d < 500$), they additionally propose to
 347 only carry out optimizations using Equation (11) while keeping single features individually (i.e. setting
 348 $|K| = 1$) and constructing synthetic data points through

$$x'_i = \begin{cases} x_i & \text{if } i \in K \\ \frac{1}{|K|+1} \cdot \left(x'_{lopt}(\emptyset) + \sum_{i \in K} x'_{lopt}(\{i\}) \right) & \text{else.} \end{cases} \quad (12)$$

349

Table 2. Mean and standard deviation of SHAP performance for ERP and MVTec data when using mean of training data (mean), zero vector (zeros), nearest neighbor in training data (NN), and optimized data points (lopt) as reference r .

(A) ERP			(B) MVTec		
r	ROC	COS	r	ROC	COS
mean	74.4 (17.1)	32.3 (25.1)	mean	64.5 (19.9)	-5.3 (2.8)
zeros	82.1 (14.2)	58.2 (16.3)	zeros	67.8 (19.9)	-2.8 (3.4)
NN	56.0 (15.1)	16.8 (38.0)	NN	66.8 (15.5)	-3.3 (3.5)
lopt	88.6 (11.2)	66.1 (20.5)	lopt	57.7 (21.4)	4.4 (8.3)

350 To demonstrate the effect of these different choices of reference values, we conduct an additional
 351 showcase using SHAP with different reference values r . Next to the mean of training data (mean), we
 352 demonstrate SHAP’s performance when using the zero vector as reference (zeros), which is another
 353 common choice in classification and regression settings. We also evaluate nearest neighbors of the normal
 354 training data (NN) as choice of reference, and integrate the approach of Takeishi and Kawahara (2020)
 355 into SHAP (lopt). For the lower dimensional ERP dataset we integrate the approach of Equation (12) into
 356 kernel-SHAP. For the larger dimensional MVTec dataset we integrate Equation (11) into partition-SHAP.
 357 Observing the results in Table 2, we notice that while zero values yield good explanations on ERP data the
 358 optimization procedure of Takeishi and Kawahara (2020) is capable of further improving results. For the
 359 image dataset MVTec, however, only minor increases on some performance metrics are observed, with the
 360 overall explanations still very poorly correlating with the ground truth.

361 Investigating the generated data points x'_{lopt} for the MVTec data in detail reveals that this approach
 362 produces many adversarial examples, i.e. examples that appear normal to the anomaly detection model,
 363 but do not truly conform to the characteristics of normal behavior. We demonstrate this behavior on
 364 our previously used anomaly x in Figure 4. Here, optimization yields a data point x'_{lopt} that is visually
 365 indistinguishable from x , with actual differences between the points enlarged in Figure 4C. This adversarial
 366 behavior indicates that the method relies on areas where the decision boundary of the underlying anomaly
 367 detection model m is not capable of generalizing and falsely associates data points with anomalous
 368 characteristics within the normal data.

369 As the generation of adversarial samples might skew the resulting explanations, future research might
 370 gain improvements over the work of Takeishi and Kawahara (2020) by specifically tuning the optimization

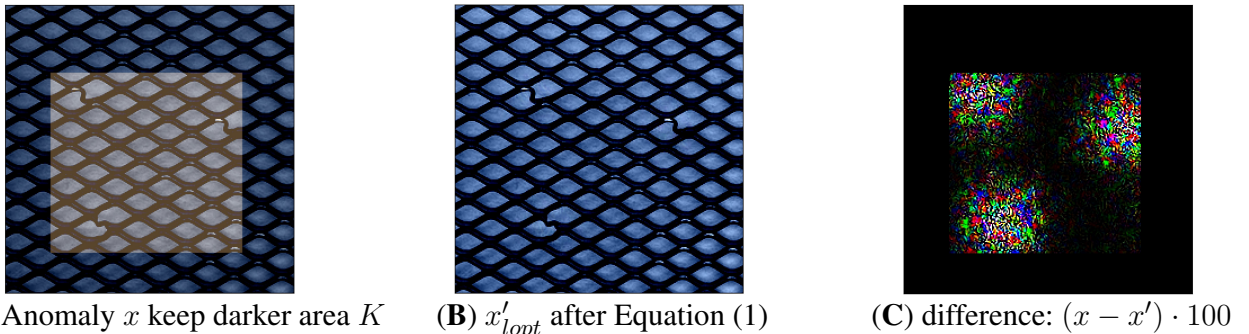


Figure 4. Perturbation with Equation (11) on MVTec: when perturbing x while keeping the dark area K shown in A, Equation (11) generates data point x'_{lopt} that is visually indistinguishable from x (B). We visualize the amplified change in pixel values in C.

371 process to find plausible inputs, which is a common technique used within the research area of counterfactual
 372 explanations (Guidotti, 2022). Additionally, the procedure of Takeishi and Kawahara (2020) introduces a
 373 large computational overhead for perturbation approaches that take thousands of sampled x' values for
 374 each data point x to explain. Further improving the performance aspects of this procedure is therefore
 375 another promising area of research.

5 GRADIENT-BASED EXPLANATIONS

376 In contrast to model-agnostic XAI approaches that base their explanations entirely on the input x and
 377 output of the investigated model $m(x)$, gradient-based approaches leverage the gradient of the model
 378 output with respect to the input $\frac{\partial m(x)}{\partial x}$ as additional information, therefore requiring investigated models
 379 to be differentiable with regards to their input and assuming that the model parameters are available during
 380 inference.

381 5.1 Saliency

382 Simonyan et al. (2014) established the use of the gradient of the output with respect to the input as a
 383 way to interpret backpropagation-based anomaly detectors. For their feature relevance explanations on
 384 image classification, which they refer to as saliency maps, they take the absolute gradient of the output
 385 with respect to the input, using the maximum gradient value for each pixel over all color channels in the
 386 case of rgb images:

$$f_{Saliency}(x, m) = \left| \frac{\partial m(x)}{\partial x} \right| \quad (13)$$

387 Beyond the utilization of the raw gradient, many applications also include a multiplication of the signed
 388 gradient values with the original input to achieve a less noisy feature relevance output (Shrikumar et al.,
 389 2016), leading to an approach commonly referred to as gradient \times input:

$$f_{gradient \times input}(x, m) = \frac{\partial m(x)}{\partial x} \cdot x \quad (14)$$

390 Nguyen et al. (2019) employ Saliency to obtain gradient-based feature relevance explanations for
 391 variational autoencoder networks on anomaly detection in NetFlow data, and further cluster the obtained
 392 feature relevance explanations to identify characteristics of anomalies.

393 5.2 Integrated Gradients

394 Sundararajan et al. (2017) note that Saliency approaches break a desirable sensitivity property that
 395 explanation approaches ought to satisfy: when only a single feature is changed within a data point, and
 396 this change alters the model's prediction, the feature should obtain a non-zero contribution. Since Saliency
 397 may violate this property in areas where the gradients are zero (e.g. around saturated activation functions),
 398 Sundararajan et al. (2017) propose a path-based approach. For a given data point x , they propose to use a
 399 reference data point r and define a smooth function giving interpolated data points on the straight-line path
 400 between x and r as $\gamma(x, r, \alpha) : \mathbb{R}^d, \mathbb{R}^d, [0, 1] \rightarrow \mathbb{R}^d$. Gradients are then calculated for these synthetic data

401 points, and an overall feature relevance explanation accumulated through a path integral

$$f_{IG}(x, m, r) = \int_{\alpha=0}^1 \frac{\partial m(\gamma(x, r, \alpha))}{\partial \gamma_i(x, r, \alpha)} \frac{\partial \gamma_i(x, r, \alpha)}{\partial \alpha} d\alpha, \quad (15)$$

402 with $\frac{\partial m(x)}{x_i}$ as the gradient of m at x along dimension i . The resulting approach, called Integrated
403 Gradients (IG), yields feature relevance values that sum to the difference of the model output at the data
404 point to be explained and the output at the reference point.

405 Sippl (2020) apply IG on anomaly detectors trained through negative sampling by choosing the nearest
406 neighbors of data points in Euclidean space as reference r . Sippl and Youssef (2022) motivate the use
407 of IG in anomaly detection from a human perspective and apply IG to real world data while sampling
408 reference points r from clustered normal data.

409 5.3 Layerwise Relevance Propagation

410 Instead of utilizing the gradient directly for feature relevance attribution, Layerwise Relevance
411 Propagation (LRP) (Bach et al., 2015) utilizes deep Taylor expansion Montavon et al. (2017) to build
412 feature relevance explanations within neural networks.

413 Consider a neural network that consists of L subsequent layers with u_i^l being the i th intermediate neuron
414 in layer $l \in \{1, 2, \dots, L - 1\}$, and where $u^1 = x$ denotes the input layer and u^L denotes the output layer.
415 LRP then computes a relevance value R_i^l for each neuron u_i^l within the network. To obtain the relevance
416 values for the input layer that correspond to feature relevance explanations $f_{LRP}(x, m) = R^1$, LRP first
417 assigns the relevance of the last network layer to the final model output ($R^L = u^L = m(x)$). Then, the
418 entire relevance is propagated to the previous layer recursively while maintaining the same total relevance
419 in each layer ($\sum_i R_i^l = \sum_j R_j^{l+1}$ for all i neurons in layer l and all j neurons in layer $l + 1$), called the
420 conservation property of LRP. The actual propagation of relevance to a neuron i of the previous layer is
421 realized through a Taylor expansion around a manually chosen root point $\tilde{u}_i^{(j)}$ with

$$R_i^l = \sum_j \frac{\partial R_j^{l+1}}{\partial u_i^l} \Big|_{\tilde{u}_i^{(j)}} \cdot (u_i - \tilde{u}_i^{(j)}). \quad (16)$$

422 While it has been shown that under specific parameter choices LRP is equivalent to the gradient \times input
423 approach in Equation (14) (Shrikumar et al., 2016), advantages of this approach are the possibility to
424 manually choose the order of Taylor expansion for each layer, which allows the approach to go beyond the
425 first order approximations of gradients when needed. Additionally, the root point \tilde{u} also needs to be chosen
426 manually for each layer, such that the conservation property of LRP is retained.

427 Amarasinghe et al. (2018) apply LRP in its standard setting on the task of detecting denial of service
428 attacks, but model the task as direct classification using feed forward neural networks instead of anomaly
429 detection architectures. As a direct application on anomaly detection architectures, Ravi et al. (2021) use a
430 standard variant of LRP that is equivalent to gradient \times input on autoencoder neural networks trained on the
431 MVTEC dataset. To appropriately adjust LRP to the task of anomaly detection, Kauffmann et al. (2020b)
432 propose specific propagation rules for common neural network layers in anomaly detection, and introduce
433 a unifying framework that transfers existing anomaly detectors into neural network representations that use

Table 3. Mean and standard deviation of gradient XAI performance comparing to ground truth explanations over all anomalies for ERP and MVTEc data respectively.

(A) ERP			(B) MVTEc		
XAI	ROC	COS	XAI	ROC	COS
noise	22.7 (7.0)	-28.7 (5.3)	noise	50.2 (1.8)	6.3 (2.7)
noise×input	52.3 (13.0)	-27.2 (8.2)	noise×input	32.3 (10.7)	-4.5 (5.5)
Saliency	50.4 (15.9)	6.0 (18.8)	Saliency	72.4 (5.0)	22.1 (8.0)
gradient×input	88.1 (13.0)	63.7 (18.5)	gradient×input	76.5 (4.4)	25.2 (8.7)
IG	78.8 (14.9)	35.6 (20.7)	IG	64.1 (6.4)	13.5 (7.4)
LRP	65.3 (20.5)	-22.0 (13.9)	LRP	65.0 (7.1)	1.8 (3.4)

434 layers for which LRP rules are defined. Through this transfer procedure, they show that LRP is applicable
 435 to a wide range of anomaly detectors.

436 5.4 Showcase and Limitations

437 5.4.1 Showcase of gradient-based approaches

438 To discuss the limitations of the introduced gradient-based approaches in detail, we again first showcase
 439 their performance in their default configuration, using the mean of training data as reference point r for IG
 440 and employing the parameter choices of Kauffmann et al. (2020b) for LRP. The resulting explanations in
 441 Table 3 compared to our random noise baselines show that all approaches are capable of finding relevant
 442 features. Especially the gradient×input approach shows strong performance on both datasets. While the
 443 multiplication with input appears necessary on the ERP data, the raw gradient of the Saliency method
 444 reaches comparable performance on the MVTEc image data. IG performs well on ERP data but struggles
 445 on the MVTEc image dataset in its default configuration, and LRP shows low performance on both datasets.

446 5.4.2 References r for path-based approaches

447 While the results of the raw gradients in the Saliency and gradient×input methods are in line with
 448 observations that the gradient signal does indeed yield explanation properties (Simonyan et al., 2014),
 449 many works in the past identified that these explanations are noisy and insensitive to specific signals (e.g.
 450 when gradients vanish due to saturated activation functions) (Shrikumar et al., 2016; Sundararajan et al.,
 451 2017). One of the proposed solutions, summing gradients along a path to avoid regions where gradients
 452 are zero as done in IG, again requires a reference data point as hyperparameter. According to the authors,
 453 this reference should be chosen to remove signal (Sundararajan et al., 2017), opening up gradient based
 454 approaches to the same issues as perturbation-based approaches with regards to finding a specific reference
 455 value that is devoid of anomaly signal, as discussed in Section 4.3.

456 To show the impact of the choice of reference r on path-based explanations, we demonstrate the effect
 457 of both established references from image classification such as the mean of training data (mean) and the
 458 zero vector (zeros), as well as the anomaly detection specific choices of nearest neighbors (NN) and the
 459 optimization scheme in Equation (12) (lopt) which we introduced in Section 4.3. While results in Table 4
 460 show decent performance of IG when using the mean and zeros references from image classification, the
 461 nearest neighbor reference performs poorly on the ERP data. The optimization scheme of Takeishi and
 462 Kawahara (2020) on the other hand indeed improves performance considerably, yielding very high XAI
 463 performance scores on all metrics for both datasets.

Table 4. Mean and standard deviation of IG performance on ERP and MVTec data with varying reference point r .

(A) ERP			(B) MVTec		
r	ROC	COS	r	ROC	COS
mean	78.8 (14.9)	35.6 (20.7)	mean	64.1 (6.4)	13.5 (7.4)
zeros	84.4 (14.0)	58.2 (20.1)	zeros	67.5 (6.3)	17.8 (8.1)
NN	54.9 (13.3)	16.5 (37.4)	NN	70.5 (5.4)	21.0 (9.1)
lopt	90.8 (11.9)	65.8 (21.7)	lopt	96.2 (2.4)	34.5 (10.5)

464 Despite the strong performance, however, an inspection of the created reference points in Figure 5 again
 465 shows that this procedure creates adversarial reference points that might skew explanations away from
 466 truly meaningful characteristics learned by the model. As seen in Figure 5B, references created through
 467 Equation (12) are visually indistinguishable from the original data point in Figure 5A and still retain their
 468 anomalous segments (previously highlighted in Figure 2D). While the changed feature values in r , which
 469 we visualized in an amplified form in Figure 5C shows that changes were indeed made in the vicinity of
 470 the three anomalous segments within anomaly x , the interpretation of explanations that result from using
 471 adversarial reference points r that contain normal behavior only for the anomaly detector but not for a
 472 human observer is unclear.

473 5.4.3 Architectural limitations of Layerwise Relevance Propagation

474 The alternative approach of LRP avoids the use of reference data points. However, the demonstrated
 475 results of Table 3 showed poor performance of LRP compared to other gradient-based approaches. Reasons
 476 for this behavior may be found in the architectural limitations of the LRP framework: while Kauffmann
 477 et al. (2020b) propose LRP rules that allow it’s application on many established differentiable anomaly
 478 detection models, the LRP framework is not capable of distributing relevance in scenarios where one
 479 layer has multiple input layers. To model common anomaly detection architectures such as autoencoder
 480 networks, where the anomaly score is usually extracted from a distance between the input layer and the
 481 reconstruction layer of the autoencoder, Kauffmann et al. (2020b) model the input layer as constant in the
 482 distance calculation. While this is a necessary assumption to retain the relevance conservation property
 483 of LRP, experimental results on the ERP autoencoder show that performance suffers significantly by not
 484 assigning a gradient to the input layer, causing LRP to generate considerably lower explanation scores

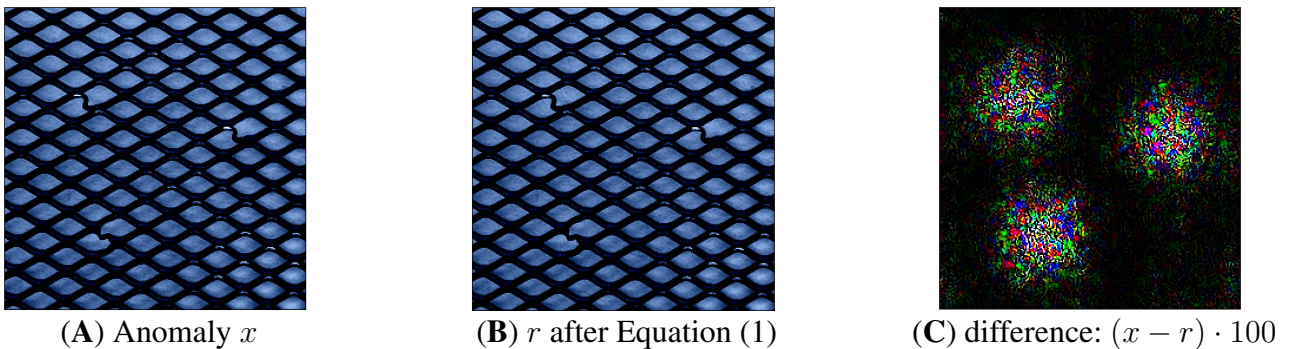
**Figure 5.** Generating reference r through Equation (11) on MVTec: similar to the perturbation issues described in Section 4.3, Equation (11) generates reference points r that are visually indistinguishable from x (A and B). We again visualize the amplified change in pixel values in C.

Table 5. Results of LRP on the ERP autoencoder when keeping the distance layer input constant as in (Kauffmann et al., 2020b) and when allowing a gradient flow. Performance improves significantly when breaking the conservation property and allowing gradient flow.

xai	variant	ROC	Cos
lrp	constant	65.3 (20.5)	-22.0 (13.9)
lrp	gradient	88.0 (12.8)	62.2 (19.2)

485 in comparison to other gradient-based approaches in Table 3A. Removing this assumption and applying
 486 the LRP variant of Kauffmann et al. (2020b) on the ERP autoencoder while retaining the gradient in the
 487 distance calculation significantly improves performance as shown in Table 5, but violates the relevance
 488 conservation property of LRP. As a result, while LRP successfully avoids the use of a reference data point, it
 489 is not readily applicable to common architectural choices such as distance calculations or skip connections.
 490 Further research into correctly distributing attribution according to the LRP properties between multiple
 491 layers that each possess a gradient with respect to the input is therefore desirable.

6 MODEL-SPECIFIC EXPLANATIONS

492 Aside from the previously introduced approaches that operate either entirely model-agnostic or only
 493 require a differentiable anomaly detector, multiple works have been proposed to generate feature relevance
 494 explanations for specific anomaly detection architectures. In contrast to the previously discussed approaches
 495 these methods heavily exploit the structure of the underlying anomaly detector to generate feature relevance
 496 explanations.

497 6.1 Depth-based Isolation Forest Feature Importance

498 Carletti et al. (2020) introduce Depth-based Isolation Forest Feature Importance (DIFFI) as an explanation
 499 approach for the well known isolation forest Liu et al. (2008) algorithm. Isolation forest is an unsupervised
 500 algorithm that uses the concept of isolation to identify anomalies using an ensemble of decision trees.
 501 The decision trees are generated by randomly splitting the training data until all training points are fully
 502 separated. Anomalies are then detected by measuring how fast they arrive on the leaf nodes of the learned
 503 trees, noting that points that are quickly isolated at random carry anomalous characteristics that allowed for
 504 the isolation. To generate feature relevance scores for single data point decisions made by isolation forests,
 505 Carletti et al. (2020) utilize this intuition by traversing a learned tree to the data point and assigning the
 506 inverse height of the data point within the tree as relevance to all features that were used as split criteria
 507 along the path to the data point. This process is repeated for all trees and feature relevance scores are
 508 summed, attributing the isolation of an individual data point to the used splitting features along all paths.
 509 Finally, all features are weighted by their inverse occurrence along all paths to counteract an effect on the
 510 explanations through the random selection during training of the isolation forest.

511 Kartha et al. (2021) extend this approach to additionally factor in the imbalance of trees before and after
 512 a split criterion, giving more relevance to features that truly isolated the data point to be explained instead
 513 of relying purely on the height of the split criterion in the tree.

514 6.2 Principal Component Analysis-based Anomaly Detection

515 Takeishi (2019) presents an approach to extract feature relevance explanations from an anomaly detector
 516 based on probabilistic principal component analysis (PCA) (Tipping and Bishop, 1999). This detector learns

517 a linear encoding $e : \mathcal{X} \rightarrow \mathcal{Z}$ of data points $\mathcal{X} \subseteq \mathbb{R}^d$ into a latent space $\mathcal{Z} \subseteq \mathbb{R}^p$ with dimensionality $p < d$
518 where the data \mathcal{X} is decomposed into its eigenvectors and only the p dimensions with highest eigenvalues
519 are retained. Points are then reconstructed through an additional linear decoding function $d : \mathcal{Z} \rightarrow \mathcal{X}$ and a
520 score of outlierness is obtained through the reconstruction error of applying the transformation through
521 $\|x - d(e(x))\|_2$ for a given data point x .

522 On this linear anomaly detector, Takeishi (2019) obtains feature relevance explanations through Shapley
523 values as described in Section 4.2. While the perturbation approaches of Section 4.2 use reference data r to
524 assess the detection output in absence of different features, Takeishi (2019) avoids the use of reference data
525 through calculating the probabilities of removed feature entries directly using the probabilistic component
526 of the anomaly detector.

527 6.3 Neuralization

528 Kauffmann et al. (2020a) introduce a "neuralization" step for explaining the outputs of one-class support
529 vector machines (OC-SVM) Schölkopf et al. (2001). In contrast to other model-specific approaches, they
530 do not explain the OC-SVM model directly but introduce a specific transfer procedure, neuralization,
531 that converts a fully trained OC-SVM into a neural network representation, allowing the subsequent
532 application of gradient-based explanation approaches such as the works discussed in Section 5. Their
533 proposed procedure transfers the final outlier scoring function learned by the OC-SVM to a two-layer
534 neural network that mimicks the behavior of the OC-SVM. Through this conversion they are able to apply
535 an LRP-style XAI approach as introduced in Section 5.3 to generate feature relevance explanations. The
536 authors further apply this "neuralization" approach to the anomaly detection approach of kernel density
537 estimation (Rosenblatt, 1956) in subsequent work (Kauffmann et al., 2020b).

538 6.4 Limitations

539 The development of highly model-specific XAI approaches bears significant potential in multiple areas.
540 While the close connection to the model architecture might allow for improved computational efficiency
541 over model-agnostic approaches (Carletti et al., 2020), the exploitation of model characteristics is also
542 a promising way to circumvent current issues of feature relevance XAI approaches such as the choice
543 of reference data as demonstrated by Takeishi (2019) on PCA. Finally, mapping fully trained anomaly
544 detection models to alternative representations as done by Kauffmann et al. (2020a) is a promising procedure
545 that allows the re-use of XAI approaches that have been identified as reliable in the domain.

546 While the continuous development of model-specific explanations approaches can therefore provide
547 numerous benefits to the domain of feature relevance XAI in anomaly detection, the main limitation of this
548 type of approach is the restriction to the specific anomaly detection model. In areas where explainability is
549 considered as a requirement of anomaly detectors, this may limit the performance of available detectors
550 in cases where a model-specific explanation framework is not available for the best performing anomaly
551 detection architecture. Especially on ERP data, hyperparameter studies of Tritscher et al. (2022b) show
552 isolation forests and PCA-based anomaly detection to perform considerably worse than other architectures,
553 which limits the application of model-specific XAI approaches such as DIFFI or the Takeishi (2019) method
554 for explaining anomaly detection of PCA. Beyond potential limitations of anomaly detection performance,
555 the promising procedure of Kauffmann et al. (2020a), who map anomaly detectors to different architectures,
556 also comes with the limitations of the XAI approach that is applied after the mapping, requiring not only
557 the mapping itself but also an XAI approach that is capable of producing reliable explanations on the
558 resulting mapped architecture.

7 DISCUSSION

559 In this work, we reviewed XAI approaches that explain single decisions of anomaly detectors by
560 highlighting which features are most anomalous. We systematically structured these feature relevance XAI
561 approaches by their access to training data and anomaly detector. We introduced the feature relevance
562 approaches and their existing adaptations to anomaly detection in detail, and showcased their current
563 limitations.

564 We showed that the many highly performing XAI approaches employed in anomaly detection require the
565 manual selection of a reference data point. This proves problematic in anomaly detection as commonly used
566 choices for reference data from other domains such as classification do not transfer to anomaly detection.

567 One approach that addresses this problem by finding optimal reference data through optimization
568 considerably improves XAI performance in our showcase, but suffers from generating adversarial data
569 points that fall outside the training data manifold. As this issue is commonly investigated within the research
570 area of counterfactual explanations (Guidotti, 2022), incorporating techniques to avoid these adversarial
571 data points during optimization constitutes a promising area for future work.

572 As another approach to circumvent issues that arise from reference data points in anomaly detection, we
573 discussed model-specific XAIs that use the model architecture to avoid the use of reference data entirely.
574 While this is a promising solution to avoid common issues with reference data, this area of research requires
575 specific design decisions for individual anomaly detectors. Therefore, developing model-specific XAI
576 approaches to ensure that state-of-the-art architectures can be explained without the use of reference data is
577 an interesting research direction.

578 Finally, once reliable XAI approaches are found within the anomaly detection domain, the extension of
579 conversion procedures that transfer trained anomaly detectors such as one-class support vector machines or
580 kernel density estimation to a more easily interpretable framework becomes a promising research area that
581 allows the transfer of reliable XAI approaches to state-of-the-art architectures.

CONFLICT OF INTEREST STATEMENT

582 The authors declare that the research was conducted in the absence of any commercial or financial
583 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

584 JT designed the showcases and wrote the paper. AK and AH critically revised the paper and experiments to
585 meet high research standards.

ACKNOWLEDGMENTS

586 We thank Daniel Schlör for insightful discussions and the anonymous reviewers for their valuable feedback
587 that helped improve this work.

DATA AVAILABILITY STATEMENT

588 The showcases contained in this study were conducted using publicly available data. Code for the
589 showcases and the used anomaly detection models are available at [https://professor-x.de/](https://professor-x.de/feature-relevance-AD)
590 [feature-relevance-AD](https://professor-x.de/feature-relevance-AD).

REFERENCES

- 591 Amarasinghe, K., Kenney, K., and Manic, M. (2018). Toward explainable deep neural network based
592 anomaly detection. In *2018 11th International Conference on Human System Interaction (HSI)* (IEEE),
593 311–317
- 594 Ancona, M., Ceolini, E., Öztireli, C., and Gross, M. (2019). Gradient-based attribution methods. In
595 *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning* (Springer). 169–191
- 596 Antwarg, L., Miller, R. M., Shapira, B., and Rokach, L. (2021). Explaining anomalies detected by
597 autoencoders using shapley additive explanations. *Expert Systems with Applications* 186, 115736
- 598 Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., et al. (2020).
599 Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward
600 responsible ai. *Information fusion* 58, 82–115
- 601 Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., and Samek, W. (2015). On pixel-wise
602 explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one* 10,
603 e0130140
- 604 Bergmann, P., Fauser, M., Sattlegger, D., and Steger, C. (2019). Mvtec ad—a comprehensive real-world
605 dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF conference on computer*
606 *vision and pattern recognition*. 9592–9600
- 607 Boser, B. E., Guyon, I. M., and Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers.
608 In *Proceedings of the fifth annual workshop on Computational learning theory*. 144–152
- 609 Carletti, M., Terzi, M., and Susto, G. A. (2020). Interpretable anomaly detection with diffi: Depth-based
610 isolation forest feature importance. *arXiv preprint arXiv:2007.11117*
- 611 Chandola, V., Banerjee, A., and Kumar, V. (2009). Anomaly detection: A survey. *ACM computing surveys*
612 *(CSUR)* 41, 1–58
- 613 Féraud, R. and Clérot, F. (2002). A methodology to explain neural network classification. *Neural networks*
614 15, 237–246
- 615 Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning* (MIT Press). [http://www.](http://www.deeplearningbook.org)
616 [deeplearningbook.org](http://www.deeplearningbook.org)
- 617 Goodman, B. and Flaxman, S. (2017). European Union Regulations on Algorithmic Decision-Making and
618 a “Right to Explanation”. *AI Magazine* 38, 50–57
- 619 Guidotti, R. (2022). Counterfactual explanations and how to find them: literature review and benchmarking.
620 *Data Mining and Knowledge Discovery* , 1–55
- 621 Hägele, M., Seegerer, P., Lapuschkin, S., Bockmayr, M., Samek, W., Klauschen, F., et al. (2020).
622 Resolving challenges in deep learning-based analyses of histopathological images using explanation
623 methods. *Scientific reports* 10, 1–12
- 624 He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In
625 *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778
- 626 Kartha, N. S., Gautrais, C., and Vercruyssen, V. (2021). Why are you weird? infusing interpretability in
627 isolation forest for anomaly detection. In *Proceedings of the Explainable Agency in AI Workshop (AAAI*
628 *2021)*. 51–57

- 629 Kauffmann, J., Müller, K.-R., and Montavon, G. (2020a). Towards explaining anomalies: A deep Taylor
630 decomposition of one-class models. *Pattern Recognit.* 101, 107198
- 631 Kauffmann, J., Ruff, L., Montavon, G., and Müller, K.-R. (2020b). The clever hans effect in anomaly
632 detection. *arXiv preprint arXiv:2006.10609*
- 633 Léger, P., Robert, J., Babin, G., Pellerin, R., and Wagner, B. (2007). Erpsim. *ERPsim Lab (erpsim. hec.
634 ca), HEC Montreal, Montreal, Qc*
- 635 Liu, F. T., Ting, K. M., and Zhou, Z.-H. (2008). Isolation forest. In *2008 Eighth IEEE Int. Conf. on Data
636 Mining (IEEE)*, 413–422
- 637 Liu, N., Shin, D., and Hu, X. (2018). Contextual outlier interpretation. In *Proceedings of the Twenty-Seventh
638 International Joint Conference on Artificial Intelligence, IJCAI-18 (International Joint Conferences on
639 Artificial Intelligence Organization)*, 2461–2467
- 640 Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *Advances
641 in Neural Information Processing Systems*. 4765–4774
- 642 Montavon, G., Lapuschkin, S., Binder, A., Samek, W., and Müller, K.-R. (2017). Explaining nonlinear
643 classification decisions with deep taylor decomposition. *Pattern recognition* 65, 211–222
- 644 Müller, R., Schreyer, M., Sattarov, T., and Borth, D. (2022). RESHAPE: Explaining Accounting Anomalies
645 in Financial Statement Audits by enhancing SHapley Additive exPlanations. In *3rd ACM International
646 Conference on AI in Finance (New York NY USA: ACM)*, 174–182
- 647 Nguyen, Q. P., Lim, K. W., Divakaran, D. M., Low, K. H., and Chan, M. C. (2019). Gee: A gradient-based
648 explainable variational autoencoder for network anomaly detection. In *2019 IEEE Conference on
649 Communications and Network Security (CNS) (IEEE)*, 91–99
- 650 Nonnenmacher, J., Holte, N.-C., and Gómez, J. M. (2022). Tell Me Why - A Systematic Literature Review
651 on Outlier Explanation for Tabular Data. In *2022 3rd International Conference on Pattern Recognition
652 and Machine Learning (PRML)*. 416–423
- 653 Owen, G. (1977). Values of games with a priori unions. In *Mathematical economics and game theory
654 (Springer)*. 76–88
- 655 Panjei, E., Gruenwald, L., Leal, E., Nguyen, C., and Silvia, S. (2022). A survey on outlier explanations.
656 *The VLDB Journal* , 1–32
- 657 Ravi, A., Yu, X., Santelices, I., Karray, F., and Fidan, B. (2021). General frameworks for anomaly detection
658 explainability: Comparative study. In *2021 IEEE International Conference on Autonomous Systems
659 (ICAS) (IEEE)*, 1–5
- 660 Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). Why should i trust you?: Explaining the predictions
661 of any classifier. In *22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data mining (ACM)*,
662 1135–1144
- 663 Robnik-Šikonja, M. and Kononenko, I. (2008). Explaining classifications for individual instances. *IEEE
664 Transactions on Knowledge and Data Engineering* 20, 589–600
- 665 Rosenblatt, M. (1956). Remarks on some nonparametric estimates of a density function. *The annals of
666 mathematical statistics* , 832–837
- 667 Ruff, L., Vandermeulen, R., Goernitz, N., Deecke, L., Siddiqui, S. A., Binder, A., et al. (2018). Deep
668 One-Class Classification. In *Proceedings of the 35th International Conference on Machine Learning
669 (PMLR)*, 4393–4402
- 670 Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., et al. (2015). ImageNet Large Scale
671 Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* 115, 211–252
- 672 Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J., and Williamson, R. C. (2001). Estimating the
673 support of a high-dimensional distribution. *Neural Computation* 13, 1443–1471

- 674 Sejr, J. H. and Schneider-Kamp, A. (2021). Explainable outlier detection: What, for Whom and Why?
675 *Machine Learning with Applications* 6, 100172
- 676 Setiono, R. and Leow, W. K. (2000). Fernn: An algorithm for fast extraction of rules from neural networks.
677 *Applied Intelligence* 12, 15–25
- 678 Shapley, L. S. (1997). A value for n-person games. *Classics in game theory* 69
- 679 Shrikumar, A., Greenside, P., Shcherbina, A., and Kundaje, A. (2016). Not just a black box: Learning
680 important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*
- 681 Simonyan, K., Vedaldi, A., and Zisserman, A. (2014). Deep inside convolutional networks: Visualising
682 image classification models and saliency maps. In *ICLR (Workshop Poster)*
- 683 Sipple, J. (2020). Interpretable, Multidimensional, Multimodal Anomaly Detection with Negative Sampling
684 for Detection of Device Failure. In *Proceedings of the 37th International Conference on Machine*
685 *Learning* (PMLR), 9016–9025
- 686 Sipple, J. and Youssef, A. (2022). A general-purpose method for applying explainable ai for anomaly
687 detection. In *International Symposium on Methodologies for Intelligent Systems* (Springer), 162–174
- 688 Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic attribution for deep networks. In *34th Int. Conf.*
689 *on Machine Learning-Volume 70* (JMLR. org), 3319–3328
- 690 Takeishi, N. (2019). Shapley values of reconstruction errors of pca for explaining anomaly detection. In
691 *2019 international conference on data mining workshops (icdmw)* (IEEE), 793–798
- 692 Takeishi, N. and Kawahara, Y. (2020). On anomaly interpretation via shapley values. *arXiv preprint*
693 *arXiv:2004.04464*
- 694 Terrell, G. R. and Scott, D. W. (1992). Variable kernel density estimation. *The Annals of Statistics* ,
695 1236–1265
- 696 Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. *Journal of the Royal*
697 *Statistical Society: Series B (Statistical Methodology)* 61, 611–622
- 698 Tjoa, E. and Guan, C. (2021). A Survey on Explainable Artificial Intelligence (XAI): Towards Medical
699 XAI. *IEEE Transactions on Neural Networks and Learning Systems* 32, 4793–4813
- 700 Tritscher, J., Gwinner, F., Schlör, D., Krause, A., and Hotho, A. (2022a). Open erp system data for
701 occupational fraud detection. *arXiv preprint arXiv:2206.04460*
- 702 Tritscher, J., Schlör, D., Gwinner, F., Krause, A., and Hotho, A. (2022b). Towards explainable occupational
703 fraud detection. In *Workshop on Mining Data for Financial Applications* (Springer)
- 704 Wang, G., Han, S., Ding, E., and Huang, D. (2021). Student-teacher feature pyramid matching for
705 unsupervised anomaly detection. In *The 32nd British Machine Vision Conference BMVC 2021*
- 706 Yepmo, V., Smits, G., and Pivert, O. (2022). Anomaly explanation: A review. *Data & Knowledge*
707 *Engineering* 137, 101946
- 708 Zhang, X., Marwah, M., Lee, I.-t., Arlitt, M., and Goldwasser, D. (2019). Ace—an anomaly contribution
709 explainer for cyber-security applications. In *2019 IEEE International Conference on Big Data (Big*
710 *Data)* (IEEE), 1991–2000