

Contextual Outlier Interpretation

Ninghao Liu
Texas A&M University
College Station, Texas
nhliu43@tamu.edu

Donghwa Shin
Texas A&M University
College Station, Texas
donghwa_shin@tamu.edu

Xia Hu
Texas A&M University
College Station, Texas
xiahu@tamu.edu

ABSTRACT

Outlier detection plays an essential role in many data-driven applications to identify isolated instances that are different from the majority. While many statistical learning and data mining techniques have been used for developing more effective outlier detection algorithms, the interpretation of detected outliers does not receive much attention. Interpretation is becoming increasingly important to help people trust and evaluate the developed models through providing intrinsic reasons why the certain outliers are chosen. It is difficult, if not impossible, to simply apply feature selection for explaining outliers due to the distinct characteristics of various detection models, complicated structures of data in certain applications, and imbalanced distribution of outliers and normal instances. In addition, the role of contrastive contexts where outliers locate, as well as the relation between outliers and contexts, are usually overlooked in interpretation. To tackle the issues above, in this paper, we propose a novel Contextual Outlier INterpretation (COIN) method to explain the abnormality of existing outliers spotted by detectors. The interpretability for an outlier is achieved from three aspects: outlierness score, attributes that contribute to the abnormality, and contextual description of its neighborhoods. Experimental results on various types of datasets demonstrate the flexibility and effectiveness of the proposed framework compared with existing interpretation approaches.

1 INTRODUCTION

Outlier detection has become a fundamental task in many data-driven applications. Outliers refer to isolated instances that do not conform to expected normal patterns in a dataset [9, 11]. Typical examples include notable human behaviors in static environment [52], online spam detection [31, 43, 55], public disease outbreaks [51], and dramatic changes in temporal signals [32, 54]. In addition, outlier detection also plays an essential role in detecting malevolence and contamination towards a secure and trustworthy cyberspace, including detecting spammers in social media [3, 53] and fraudsters in financial systems [38].

Complementing existing work on detecting outliers, interpretability of the detection results is becoming increasingly important for domain experts especially those with limited data science background [25]. First, complicated statistical inferences and algorithms impede the domain experts from understanding and trusting the outlier detection methods. The focus of existing techniques is to efficiently and effectively detect outliers by tackling the challenges including the curse of dimensionality [2, 16, 24], the massive data volume [4, 40], and heterogeneous information sources [17, 36].

However, the essential reasons that cause the abnormality of outliers are usually ignored and cannot be revealed explicitly with the detection outcome to end users. Second, it is difficult for end users to comprehensively evaluate the outlier detection performance. It is time-consuming and labor-intensive to manually examine the detection results without an intuitive understanding of the outliers. Current evaluation metrics such as area under ROC curve (AUC) and nDCG [9] only provide limited information about the intuitive characteristics of the outliers. Also, a detection method that works effectively in one dataset or application is not guaranteed to have good performance in others. Unlike supervised learning methods, outlier detection is usually implemented with unsupervised methods and cannot be evaluated in the same way. Thus, effective outlier interpretation would significantly facilitate the usability of different types of outlier detection methods in real-world applications.

To this end, one straightforward way for outlier interpretation is to apply feature selection to identify a subset of original attributes that distinguish outliers from normal instances [13, 21, 33, 50]. However, first it is difficult for some existing methods to efficiently handle datasets of large size or high dimensions [50], or effectively obtain interpretations from complex data types and distributions [13, 21]. Second, we also want to measure the abnormality level of each outlier through interpretation process. Outliers have different levels of abnormality. The results provided by detectors could be binary labels indicating whether each data instance is an outlier or not. Even if abnormality scores are estimated along with data instances, they are usually in different scales when different detection methods are applied. A unified scoring formula provided through interpretation will facilitate the comparisons among various detectors. Third, besides focusing on discovering notable attributes of outliers, we would also like to analyze the context (e.g., contrastive neighborhood) in which outliers are detected. "It takes two to tango." Discovering the relations between an outlier and its context for contrast would provide richer information before taking actions to deal with the outlying objects in real applications.

To tackle the aforementioned challenges, in this paper, we propose a novel Contextual Outlier INterpretation (COIN) approach to provide explanations for outliers identified by detectors. We define the interpretation of an outlier as the triple of noteworthy features, the degrees outlierness and the contrastive context with respect to the outlier query. The first two components are extracted from the relations between the outlier and its context. Also, the interpretations of all outliers can be integrated for evaluating the given outlier detection model. The performance of different detectors can also be compared through interpretations as COIN provides a unified evaluation basis. COIN can also be applied to existing outlier/anomaly detection methods which already provide explanations for their results. In addition, prior knowledge of attribute characteristics about certain application scenarios can be easily

Notation	Definition
N	the number of data points in the dataset
M	the number of attributes
\mathbf{x}	a data instance, $\mathbf{x} \in \mathbb{R}^M$
a_m	the m -th attribute
\mathcal{X}	all data instances, $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
h	an outlier detection method
\mathcal{O}	the collection of detected outliers
\mathbf{o}_i	outlier i identified by the detector
\mathcal{O}_i	the outlier class corresponding to \mathbf{o}_i
\mathcal{C}_i	the context of outlier \mathbf{o}_i
k	the number instances included in \mathcal{C}_i
$s(a_m)$	suspicious score of attribute a_m
$d(\mathbf{x})$	outlierness score of \mathbf{x}

Table 1: Symbols and Notations

incorporated into the interpretation process, which enables end users to perform model selection according to specific demands. The contributions of this work are summarized as follows:

- We design a novel interpretation approach called COIN to explain outliers. The approach is model-agnostic. Each interpretation is composed of three aspects including the outlying attributes, the outlierness score, and descriptions of the local context.
- We show that the performance of outlier detectors can be evaluated by interpreting their detection results through COIN.
- Comprehensive evaluations on interpretation quality and model selection accuracy are conducted through experiments with both real-world and simulated datasets. Case studies are also presented and discussed for intuitive explanations.

2 PRELIMINARIES

Background Many approaches have been proposed for outlier detection. These approaches can be divided into three categories: density-based, distance-based and model-based. Density-based approaches try to estimate the data distribution, where instances that fall into low-densities regions are returned as outliers [2, 8, 46]. Distance-based methods identify outliers as instances isolated far away from their neighbors [4, 7, 22, 29, 40]. For model-based ones, usually a specific model (e.g., classification, clustering or graphical model) is applied to the data, and outliers are those who do not fit the model well [19, 42, 48]. Other main focuses of relevant research include tackling the challenges of the curse of dimensionality [2, 16, 24], the massive data volumn [4, 40] and heterogeneous information sources [36]. However, interpretation of detection results is usually overlooked. Although some recent anomaly detection methods provide explanation with their outcome [17, 28, 30, 37], they do not represent all the scenarios. The ignorance of outlier interpretation may lead to several problems. First, for security-related domains, where outlier detection is widely applied, explanations affect whether the results will be accepted by end users. Second, the sparsity of outliers brings uncertainty to evaluation methods. Small disturbance on the detection results may lead to significant variations in evaluation results using traditional metrics such as ROC scores [12]. Third, it is usually difficult to obtain labels of outliers, so we wonder if it is possible to evaluate the detection

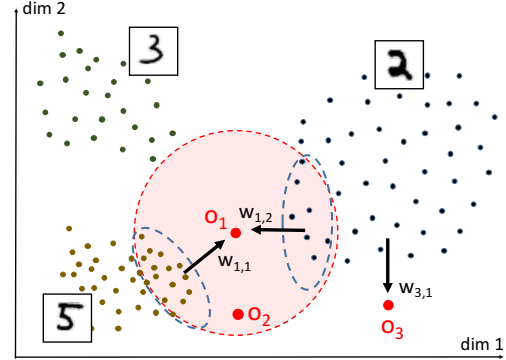


Figure 1: A toy example of outlier interpretation by resolving its context in to clusters.

performance without ground-truth labels. In this work, we resort to interpretation methods to tackle the challenges above.

Notations The notations used in this paper are introduced as below and in Table 1. Let \mathcal{X} denotes the collection of all data. N is the number of data instances in \mathcal{X} . Each data instance is denoted as $\mathbf{x} \in \mathbb{R}^M$, where M is the number of attributes. The m^{th} attribute is denoted as a_m . We use h to represent an outlier detector. The collection of outliers identified by a detector is represented as \mathcal{O} , in which a single outlier is denoted as $\mathbf{o} \in \mathbb{R}^M$. The *context* of an outlier \mathbf{o} , i.e., \mathcal{C}_i , is composed of its k -nearest normal instances. Each \mathcal{C}_i could consist of some smaller clusters $\mathcal{C}_{i,1}, \mathcal{C}_{i,2}, \dots, \mathcal{C}_{i,L}$. Among the detected outliers, some are far away from the bulks of the dataset, while others are just marginally abnormal. We define the degree of outlierness for an instance \mathbf{x} as *outlierness* denoted as $d(\mathbf{x}) \in \mathbb{R}_{\geq 0}$. The reason for clustering the context is illustrated in Figure 1. There are three clusters, each of which represents images of a digit. Red points are outliers detected by a certain algorithm. Clusters of digit “2” and “5” compose the context of outlier \mathbf{o}_1 . The interpretation of \mathbf{o}_1 , denoted as $\mathbf{w}_{1,1}$ and $\mathbf{w}_{1,2}$, can be obtained by contrasting it with the two clusters respectively. However, it would difficult to explain the outlierness of \mathbf{o}_1 if clusters of digit “2” and “5” are not differentiated.

Problem Definition Based on the analysis above, here we formally define the outlier interpretation problem as follows. Given a dataset \mathcal{X} and the query outliers \mathcal{O} detected therefrom, the *interpretation* for each outlier $\mathbf{o}_i \in \mathcal{O}$ is defined as a composite set: $\mathcal{E}_i = \{\mathcal{A}_i, d(\mathbf{o}_i), \mathcal{C}_i = \{\mathcal{C}_{i,l} | l \in [1, L]\}\}$. Here \mathcal{A}_i include the abnormal attributes of \mathbf{o}_i with respect to \mathcal{C}_i , $d(\mathbf{o}_i)$ is the outlierness score of \mathbf{o}_i , \mathcal{C}_i denotes the context of \mathbf{o}_i and $\mathcal{C}_{i,l}$ is the l -th cluster.

3 CONTEXTUAL OUTLIER INTERPRETATION FRAMEWORK

The general framework of Contextual Outlier Interpretation (COIN) is illustrated in Figure 2. Given a dataset \mathcal{X} and detected outliers \mathcal{O} , we first map the interpretation task to a classification problem due to their similar natures. Second, the classification problem over the whole data is partitioned to a series of regional problems focused on the context of each outlier query. Third, a collection of simple and local interpreters g are built around the outlier. At last, the outlying attributes and outlierness score of each outlier

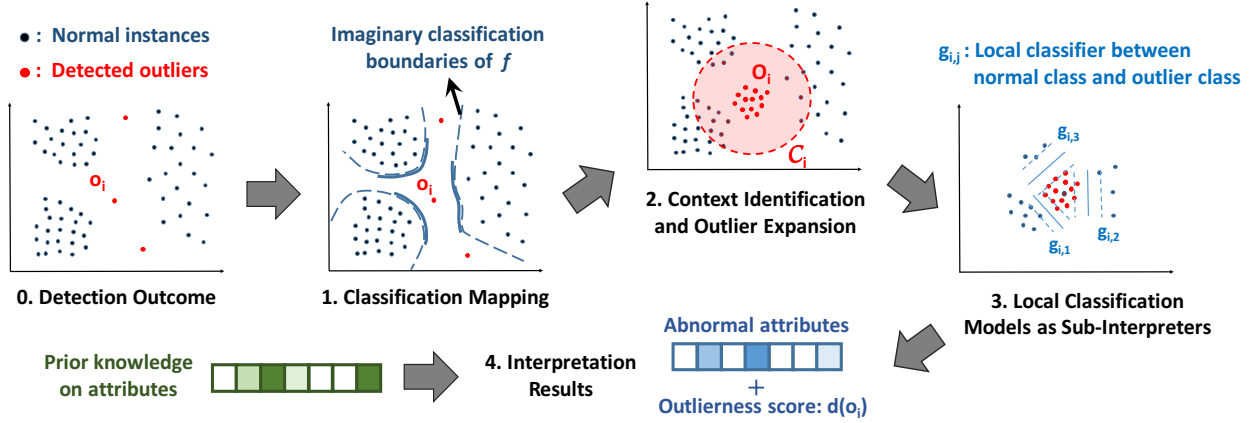


Figure 2: Contextual Outlier Interpretation Framework

can be directly obtained from the parameters of g , by combining the application-related prior knowledge. The details of each step are discussed in the following subsections.

3.1 Explain Outlier Detector with Classifiers

In this module, we will establish the correlation between outlier detection and classification. The close relationship between the two types of problems motivates us to design the interpretation framework from the classification perspective.

Formally, an outlier detection function can be denoted as $h(\mathbf{x}|\theta, \mathcal{X})$, where $\mathbf{x} \in \mathcal{X}$, θ and \mathcal{X} represent the function parameters. Here the dataset \mathcal{X} is treated as parameters since data instances affect the degree of normality of each other. The abnormality of an instance is typically represented by either a binary label or a continuous score. In the former case, an instance is categorized as either normal or abnormal, while the latter expresses the degree to which an instance is abnormal. The latter case can be easily transformed to the former if a threshold is set as the separating mark between inlier and outlier classes [2, 17, 29]. This form of binary detection motivates us to analyze the mechanism of outlier detectors using classification models. Although outlier detection is usually tackled as an unsupervised learning problem, we can assume there exists a latent hyperplane specified by certain function $f(\mathbf{x}|\theta') : \mathbb{R}^M \rightarrow \{0, 1\}$ that separates outliers from normal instances. Here θ' represents the parameters of f . This connection between outlier detection and supervised learning has also been implied in some previous work [1, 42]. An intuitive example can be found in Step 1 of Figure 2. Blue points represent normal instances, and red points indicate the detected outliers. The decision boundaries defined by f are shown using dotted curves. In this setting, the outlier detector is actually trying to mimic the behavior of the decision function.

Given the outliers \mathcal{O} identified by detector h , we want to recover the implicit decision function f which leads to the similar detection results as h . The problem is thus formulated as below,

$$\arg \min_f \mathcal{L}(h, f; \mathcal{O}, \mathcal{X} - \mathcal{O}), \quad (1)$$

where \mathcal{L} is the loss function that includes all the factors (e.g., classification error and simplicity of f) we would like to consider. \mathcal{O} and

$\mathcal{X} - \mathcal{O}$ represent outlier class and inlier class, respectively. However, the final form of f could be very complicated if outliers have diverse abnormal patterns and the whole dataset contains complex cluster structures. Such complexity prevents f from directly providing intuitive explanations for the detector. This is also a common issue in many supervised learning tasks, where the highly complicated prediction function makes the classification model almost a black box. A straightforward solution is to first obtain f and then interpret it [6, 41]. This pipeline, however, will introduce new errors in the intermediate steps, and it is more computationally expensive to deal with large datasets. An approach for directly interpreting outlier detectors is needed.

3.2 Local Interpretation for Individual Outliers

By utilizing the isolation property of outliers, we can decompose the overall problem of detector interpretation into multiple regional tasks of explaining individual outliers:

$$\begin{aligned} \min_f \mathcal{L}(h, f; \mathcal{O}, \mathcal{X} - \mathcal{O}) &\Rightarrow \min_f \sum_i \mathcal{L}(h, f; \mathbf{o}_i, \mathcal{C}_i) \\ &\Rightarrow \sum_{i \in [1, |\mathcal{O}|]} \min_{g_i} \mathcal{L}(h, g_i; \mathbf{o}_i, \mathcal{C}_i) \\ &\Rightarrow \sum_{i \in [1, |\mathcal{O}|]} \min_{g_i} \mathcal{L}(h, g_i; \mathcal{O}_i, \mathcal{C}_i). \end{aligned} \quad (2)$$

In this way, the original problem is transformed to explaining each outlier \mathbf{o}_i with respect to its context counterpart \mathcal{C}_i . Since the number of outliers is usually small, we avoid dealing with the whole dataset which could be large. Here g_i represents the local parts of f exclusively for classifying \mathbf{o}_i and \mathcal{C}_i . In Figure 2, for example, g_i is highlighted by the bold boundaries around \mathbf{o}_1 in Step 1, and \mathcal{C}_i consists of the normal instances enclosed in the circle in Step 2. Since there is a data imbalance between the two classes, by applying strategies such as synthetic sampling [18], \mathbf{o}_i is expanded to a hypothetical outlier class \mathcal{O}_i with comparable size to \mathcal{C}_i . As it is common for outlier detectors to measure the outlierness of instances based on their contexts, a proper interpretation method would better take this into consideration.

3.3 Resolve Context for Outlier Explanations

Now we focus on interpreting each single outlier \mathbf{o}_i by solving g_i in Equation 2. Since g_i is the local classifier separating \mathcal{O}_i from \mathcal{C}_i , the current task is turned into interpreting the classification boundary of g_i . Let $p_{\mathcal{O}_i}(\mathbf{x})$ and $p_{\mathcal{C}_i}(\mathbf{x})$ denote the probability density function for the outlier class and inlier class, respectively. Since the context \mathcal{C}_i for different i could have various cluster structures as shown in Figure 1, it is difficult to directly measure the degree of separation between \mathcal{O}_i and \mathcal{C}_i or to discover the attributes that characterize the differences between the two classes. Therefore, we further resolve $\mathcal{L}(h, g_i; \mathcal{O}_i, \mathcal{C}_i)$ to a set of simpler problems. According to Bayesian decision theory, the classification error equals to

$$\begin{aligned} & P^{err}(\mathcal{O}_i, \mathcal{C}_i) \quad (3) \\ &= P(\mathcal{O}_i) \int_{\mathcal{C}_i} p(\mathbf{x}|\mathcal{O}_i) d\mathbf{x} + P(\mathcal{C}_i) \int_{\mathcal{O}_i} p(\mathbf{x}|\mathcal{C}_i) d\mathbf{x} \\ &\approx \left(\sum_{l \in [1, L]} P(\mathcal{O}_i) \int_{\mathcal{C}_{i,l}} p(\mathbf{x}|\mathcal{O}_i) d\mathbf{x} \right) + \left(\sum_{l \in [1, L]} P(\mathcal{C}_{i,l}) \int_{\mathcal{O}_i} p(\mathbf{x}|\mathcal{C}_{i,l}) d\mathbf{x} \right) \\ &= \sum_{l \in [1, L]} \left(P(\mathcal{O}_i) \int_{\mathcal{C}_{i,l}} p(\mathbf{x}|\mathcal{O}_i) d\mathbf{x} + P(\mathcal{C}_{i,l}) \int_{\mathcal{O}_i} p(\mathbf{x}|\mathcal{C}_{i,l}) d\mathbf{x} \right) \\ &= \sum_{l \in [1, L]} P^{err}(\mathcal{O}_{i,l}, \mathcal{C}_{i,l}). \end{aligned}$$

Suppose we can split the context \mathcal{C}_i into multiple clusters $\{\mathcal{C}_{i,l} | l \in [1, L]\}$ that are sufficiently separated from each other, then cluster $\mathcal{C}_{i,l}$ is the only dominant class near the decision boundary between \mathcal{O}_i and $\mathcal{C}_{i,l}$. Then each term in the summation can be treated as an individual sub-problem of classification without mutual inference. By combining Equation 2 and Equation 3, our interpretation tasks is finally formulated as:

$$\min_f \mathcal{L}(h, f; \mathcal{O}, \mathcal{X} - \mathcal{O}) \Rightarrow \min_{g_{i,l}} \sum_i \sum_l \mathcal{L}(h, g_{i,l}; \mathcal{O}_{i,l}, \mathcal{C}_{i,l}). \quad (4)$$

By now we are able to classify $\mathcal{O}_{i,l}$ and $\mathcal{C}_{i,l}$ with a simple and explainable model $g_{i,l}$ such as linear models and decision trees, where the outlying attributes $\mathcal{A}_{i,l}$ can be extracted from *model parameters* [26, 41]. The overall interpretation for \mathbf{o}_i can be obtained by integrating the results across all context clusters $\mathcal{C}_{i,l}$, $l \in [1, L]$.

The estimated time complexity for implementing the framework above is $O(|\mathcal{O}| \times L \times T_g)$, where T_g is the average time cost of constructing $g_{i,l}$. Due to the scarcity of outliers, $|\mathcal{O}|$ is expected to be small. Each $g_{i,l}$ involves $\mathcal{O}_{i,l}$ and $\mathcal{C}_{i,l}$. Since $\mathcal{C}_{i,l}$ is only a small subset of data points around an outlier, and $\mathcal{O}_{i,l}$ has comparable size with $\mathcal{C}_{i,l}$, both of their cardinalities should be small, which significantly reduces the time T_g . Moreover, the interpretation processes of different outliers are independent of each other, thus can be implemented in parallel to further reduce the time cost.

4 OUTLIERNES-COUPLED SUSPICIOUS ATTRIBUTES DISCOVERY

After introducing the general framework of COIN, we have resolved the vague problem of outlier interpretation into a collection of classification tasks around individual outliers. In this section, we will propose concrete solutions for explaining an individual outlier, including discovering its *abnormal attributes* and measuring the *outlierness score*.

4.1 Context Identification and Clustering

Given an outlier \mathbf{o}_i spotted by detector h , first we need to identify its context \mathcal{C}_i in the data space. As introduced before, \mathcal{C}_i consists of the nearest neighbors of \mathbf{o}_i . Here we use Euclidean distance in attribute space as the point-to-point distance measure. The neighbors are chosen only from normal instances, because outlier instances do not represent the common patterns in the data. These nearest neighbors are regarded as the representatives for the local background around the outlier. Although these adjacent data instances are only the ‘‘tips of icebergs’’ to the whole data distribution, they are the gates of inlier regions facing outliers and thus are adequate for discriminating the two classes. An example of context identification can be found in Figure 1, in which the context instances of \mathbf{o}_3 are embraced in the red circle.

As local context may indicate some interesting structures (e.g., instances with similar semantics are located close to each other in the attribute space), it is necessary to further segment the neighbors into multiple disjoint clusters, where each cluster corresponds to one aspect of the context. Such an idea of context clustering is inspired by various anomaly detection models which perform data clustering prior to recognizing anomalies [17, 30, 39, 45]. To determine the number of clusters L in \mathcal{C}_i , we adopt the measure of *prediction strength* [47] which shows good performance even when dealing with high-dimensional data. After obtaining L , common clustering algorithms such as K-Means or hierarchical clustering methods can be applied to divide \mathcal{C}_i into multiple clusters, i.e. $\mathcal{C}_i = \{\mathcal{C}_{i,1}, \mathcal{C}_{i,2}, \dots, \mathcal{C}_{i,L}\}$. Minor clusters whose size is too small will be ignored in subsequent procedures (e.g., the data points of cluster 3 can be ignored in the context of \mathbf{o}_1 in Figure 1). This is because minor clusters are usually farther from the outlier than other major clusters, or simply represent noise in data.

4.2 Maximal-Margin Linear Explanations

Given an outlier \mathbf{o}_i and one of its context clusters $\mathcal{C}_{i,l}$, we now focus on the problem of $g_{i,l}$ solved by minimizing $P^{err}(\mathcal{O}_{i,l}, \mathcal{C}_{i,l})$. As the exact distribution formulas for the two classes are not known, we use non-parametric estimation method to model their probability distributions. We choose Parzen Windows of Gaussian distribution with diagonal covariance matrix as kernels. For a certain class \mathcal{C} , its density distribution $p(\mathbf{x}|\mathcal{C}) = \sum_{\mathbf{x}_n \in \mathcal{C}} \mathcal{N}(\mathbf{x} | \mathbf{x}_n, \sigma^2 I) / |\mathcal{C}|$, where \mathcal{N} denotes the Gaussian distribution. After plugging the expression above into $P^{err}(\mathcal{O}_{i,l}, \mathcal{C}_{i,l})$, a careful analysis [49] indicates that, the hyperplane characterized by the optimal classifier converges to the *maximal-margin* hyperplane if we set σ to be small. Then the estimated Bayes error $P^{err}(\mathcal{O}_{i,l}, \mathcal{C}_{i,l})$ is dominated by an expression proportional to the exponential in $-\text{margin}^2(g_{i,l})/(\sigma^2)$. It means that, let $d(\mathbf{o}_i, \mathcal{C}_{i,l})$ denote the separability from \mathbf{o}_i to $\mathcal{C}_{i,l}$, $d(\mathbf{o}_i, \mathcal{C}_{i,l})$ increases monotonically as the margin of the hyperplane increases. In another words, the margin of the hyperplane characterized by classifier $g_{i,l}$ can reflect the relative distance between an outlier class $\mathcal{O}_{i,l}$ and its contextual cluster $\mathcal{C}_{i,l}$.

There are several concerns with respect to choosing a concrete form of $g_{i,l}$. First, $g \in G$ should be simple to understand by end users. For examples, we may expect the number of non-zero weights to be small for linear models, or the rules to be concise in decision trees [41]. Second, since outliers are usually highly separated from

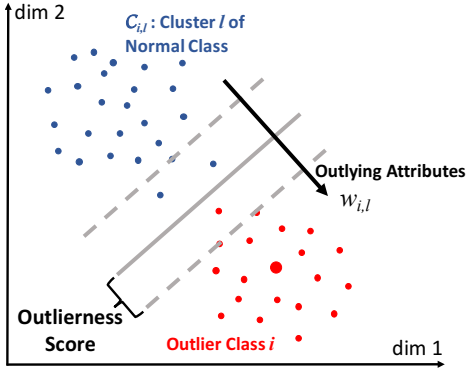


Figure 3: Local Outlier Interpretation from SVM Parameters

their context, there could be multiple solutions all of which could classify the outliers and inliers almost perfectly, so how to choose the one that best fits the mechanism which causes outliers to be susceptible? Here we let $g \in G$ belongs to linear models, i.e., $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. We impose the l_1 norm constraint on \mathbf{w} , where attributes a_m corresponding to nonzero $|w[m]|$ are reported as abnormal [14]. Motivated by the *isolation* property of outliers [29, 33], we use l_1 norm support vector machine (1-norm SVM) [56] to build g . As shown in Figure 3, outlying attributes can be identified from weight vector \mathbf{w} and the outlierness score is relevant to the margin of SVM. The local loss $\mathcal{L}(h, g_{i,l}; \mathcal{O}_{i,l}, \mathcal{C}_{i,l})$ to be minimized is as below:

$$\begin{aligned} & \sum_{n=1}^{N_{i,l}} (1 - y_n g(\mathbf{x}_n) - \xi_n)_+ + c \sum_{n=1}^{N_{i,l}} \xi_n, \\ \text{s.t. } & \xi_n \geq 0 \\ & \|\mathbf{w}\|_1 \leq b \end{aligned} \quad (5)$$

where $N_{i,l} = |\mathcal{O}_{i,l} \cup \mathcal{C}_{i,l}|$, $(\cdot)_+$ is the hinge loss, ξ_n is slack variable as we allow some instances to fall into the margin, b and c are the tuning parameters. Here $y_n = 1$ if $\mathbf{x}_n \in \mathcal{C}_{i,l}$ and $y_n = -1$ if $\mathbf{x}_n \in \mathcal{O}_{i,l}$.

From the parameters of the local model $g_{i,l}$, we are able to select the most significant attributes that make \mathbf{o}_i isolated from $\mathcal{C}_{i,l}$. In this way, we avoid searching through the exponentially large space of all possible attribute subsets. Let $\mathbf{w}_{i,l}$ denote the weight vector of $g_{i,l}$, the significance score of attribute a_m thus equals to $s_{i,l}(a_m) = |w_{i,l}[m]| / \gamma_{i,l}^m$. Here $\gamma_{i,l}^m$ denotes the resolution of attribute a_m in $\mathcal{C}_{i,l}$, i.e., the average distance between an instance and its closest neighbors in $\mathcal{C}_{i,l}$ along the m^{th} axis. The score above can be seen as the absolute value of weight $w_{i,l}[m]$ normalized by the scale of attribute m on the context cluster $\mathcal{C}_{i,l}$. Although \mathcal{X} may have been normalized before fed into the interpreter, it is still necessary to reconsider the scale of attributes in each contextual cluster, because the density of the data could vary in different localities with respect to different subsets of features [2, 34]. For discrete attributes, we may need to set a low bound on the denominator in case all neighbor instances aggregate on a single point. The attributes with large $s_{i,l}(a_m)$ constitute the set of abnormal attributes $\mathcal{A}_{i,l}$ with respect

to $\mathcal{C}_{i,l}$. The overall significance score of attribute a_m for \mathbf{o}_i is

$$s_i(a_m) = \frac{1}{|\mathcal{C}_i|} \sum_l |\mathcal{C}_{i,l}| s_{i,l}(a_m), \quad (6)$$

which is the weighted average score for a_m over all clusters, weighted by the relative size of each cluster. Attributes with large s_i scores constitute the abnormal attributes for \mathbf{o}_i .

After obtaining the local classifier $g_{i,l}$, we are able to measure the outlierness score $d(\mathbf{o}_i)$ of \mathbf{o}_i . Besides non-negativeness and finiteness, an important requirement for an outlierness measure is *ranking-stability* [23]. It is expected that $d(\mathbf{x})$ would reflect the relative degree to which \mathbf{x} deviates from its context. From the analysis in Section 4.2, we can use the *margin* of the hyperplane defined by $g_{i,l}$ as the outlierness measure of an outlier \mathbf{o}_i with respect to its normal instances counterpart $\mathcal{C}_{i,l}$, i.e., $d_l(\mathbf{o}_i) = |g_{i,l}(\mathbf{o}_i)| / \|\mathbf{w}_{i,l}\|$ where $\|\cdot\|$ is l_2 norm. This measure is robust to the high dimensionality of data, as \mathbf{w} is sparse and $d_l(\mathbf{o}_i)$ is calculated in a low dimensional space.

4.3 Incorporate Prior Knowledge into Interpretation

In real-world applications, the importance of different attributes varies according to different scenarios [10, 35, 53]. Take Twitter spammer detection as an example. We discuss two attributes of users: the number of followers (N_{fer}) and the ratio of tweets posted by API (R_{API}). A spammer tends to have small N_{fer} value as they are socially inactive, but large R_{API} in order to conveniently generate malevolent content. However, it is easy for spammers to intentionally increase their N_{fer} by following each other, while manually decreasing R_{API} is more difficult due to the expense human labor. In this sense, R_{API} is more robust and more important than N_{fer} in translating detected outliers as social spammers. To represent the different roles of attributes, we introduce two vectors $\boldsymbol{\beta}$ and \mathbf{p} , where $\beta_m \in \mathbb{R}_{\geq 0}$ denotes the relative degree of significance assigned to attribute a_m , and $p_m \in \{-1, 0, 1\}$ denotes the prior knowledge on the expected magnitude of attribute values of outliers. $p_m = -1$ means we expect outliers to have small value for a_m (e.g., N_{fer}), $p_m = 1$ means the opposite (e.g., R_{API}), while $p_m = 0$ means there is no preference. Therefore, the outlierness score of \mathbf{o}_i with respect to $\mathcal{C}_{i,l}$ is refined as:

$$d_l(\mathbf{o}_i) = \frac{|g_{i,l}(\mathbf{o}_i)|}{\gamma_{i,l} \|\mathbf{w}_{i,l}\|} \frac{\mathbf{w}'_{i,l}}{\|\mathbf{w}_{i,l}\|} \circ \boldsymbol{\beta}, \quad (7)$$

where the operator \circ denotes element-wise multiplication, $w'[m] = \min(0, w[m])$ if $p_m = 1$, and $w'[m] = \max(0, w[m])$ if $p_m = -1$. If we label outliers with 1 and inliers with -1 , the sign is reversed. The motivation of introducing \mathbf{w}' is that, if interpretation results (e.g., R_{API} is small) does not conform with the expectation expressed by the prior knowledge (e.g., R_{API} is expected to be large to signify spammers), then the outlierness score of the outlier should be deducted. Here $\gamma_{i,l}$ is the average distance from an instance to its closest neighbor in $\mathcal{C}_{i,l}$. It normalizes the outlierness measure with respect to the data density of different clusters. Therefore, the overall outlierness score for \mathbf{o}_i across all context clusters is:

$$d(\mathbf{o}_i) = \frac{1}{|\mathcal{C}_i|} \sum_l |\mathcal{C}_{i,l}| d_l(\mathbf{o}_i), \quad (8)$$

which comprehensively considers the isolation of \mathbf{o}_i over different contexts. Now we have obtained all of the three aspects of interpretation $\mathcal{E}_i = \{\mathcal{A}_i, d(\mathbf{o}_i), \mathcal{C}_i = \{\mathcal{C}_{i,l} | l \in [1, L]\}\}$. If a normal instance is misdetected as an outlier by a detection method, then \mathcal{E}_i is able to identify such mistake, since $d(\mathbf{o}_i)$ will be small and \mathcal{A}_i is less likely to conform to the prior knowledge.

5 EXPERIMENTS

In this section, we present evaluation results to assess the effectiveness of our framework. We try to answer the following questions: 1) How accurate the proposed framework is to identify the outlying attributes from outlier queries? 2) Can we faithfully measure the outlierness score of outliers? 3) How effective is the prior knowledge of attributes in refining outlier detection results? 4) Can our framework correctly evaluate the performance of outlier detectors by only using interpretation results rather than the ground truth?

5.1 Datasets

The real-world datasets used in our experiments include Wisconsin Breast Cancer (WBC) dataset [5], MNIST dataset and Twitter spammer dataset [53]. The outlier labels are available. WBC dataset records the measurements for breast cancer cases with two classes, i.e. benign and malignant. The former is considered as normal, while we downsampled 25 malignant cases as the outliers. MNIST dataset includes a collection of 28×28 images of handwritten digits. In our experiments, we only use the training set which contains 42,000 examples. Instead of using raw pixels as attributes, we build a Restricted Boltzmann Machine (RBM) with 150 latent units to map images to a higher-level attribute space [20]. The new low-dimensional attributes are more proper for interpretation than raw pixels. A multi-label logistic classifier is then built to classify different written digits, and the ground truth outliers are selected as the misclassified instances downsampled to 1,000 of them. The Twitter dataset contains information of normal users and spammers crawled from Twitter. Attributes are classified into two categories according to whether they are robust to the disguise of spammers. Low robustness attributes refer to those which can be easily controlled by spammers to avoid being detected, while high robustness attributes are more trustworthy in discriminating spammers from normal users [53].

We also build two synthetic datasets with ground truth outlying attributes for each outlier. Both datasets consist of multiple clusters as normal instances generated under multivariate Gaussian distributions. Outliers are created by distorting some samples' attribute values beyond certain clusters, while keeping other attributes within the range of the normal instances. In the first dataset, each outlier is close to only one normal cluster and far away from the others. In the second dataset, an outlier is in the vicinity of several normal clusters simultaneously, while its outlying attributes differ with respect to different neighbors, so that a more refined interpretation approach is required.

5.2 Baseline Methods

We compare COIN with some baseline methods including outlying-aspect mining techniques and classifier interpretation approaches summarized as below:

	SYN1	SYN2	WBC	Twitter	MNIST
N	405	405	458	11,000	42,000
M	15	15	9	16	150
$ \mathcal{O} $	30	30	25	1,000	1,000

Table 2: Details of the datasets in experiments

- CA-lasso (CAL) [33]: Measure the separability between outlier and inliers as the classification accuracy between the two classes, and then apply feature selection methods (e.g., LASSO) to determine the attribute subspace as explanations.
- Isolation Path Score with Beam Search (IPS-BS) [50]: Apply isolation path score [29] to measure outlierness. The score is then used to guide the search of subspaces, where Beam Search is applied as the main strategy.
- LIME [41]: An effective global classification model is first constructed to classify outliers and inliers. Then the outlying attributes for each outlier is identified by locally interpreting the classification model around the outlier. Oversampling is applied to prevent data imbalance. A neural network is used as the global classifier for MNIST data, and SVMs with RBF kernel are used for other datasets.

5.3 Outlying Attributes Evaluation

The goal of this experiment is to verify that the attributes identified by COIN are indeed outlying. Since ground-truth outlying attributes of real-world datasets are not available, we append M noise attributes to all real-world data instances. We simply assume that all of the original attributes are outlying attributes, and noise attributes are not. For each outlier, we apply our approach as well as baseline methods to infer the outlying attributes, and compare the results with the ground truth to evaluate their performances. In our experiments, we choose 8% of nearest neighbors of an outlier \mathbf{o}_i as its context \mathcal{C}_i . The radius of synthetic sampling for building the outlier class \mathcal{O}_i is set as half of the distance to the inlier class \mathcal{C}_i , in order to suppress the overlap between the two classes. The hyperparameters in SVM models are determined through validation, where some samples from \mathcal{O}_i and \mathcal{C}_i are randomly selected as the validation set. The same hyperparameter values are used for all outliers in the same dataset. We report the Precision, Recall and F1 score averaged over all the outliers queries in Table 3. Besides finding that COIN consistently indicates good performance, some observations can be made as follows:

- In general the Recall value for SYN2 is lower than that for SYN1, while the Precision value is on the contrary. This is because each outlier in SYN2 has more than one context clusters, and the real outlying attributes for each outlier vary with respect to different clusters. In this case, extracting as many ground truth attributes as possible becomes a more challenging task.
- IPS-BS is relatively cautious in making decisions. It first detects trivial abnormal attributes and will stop early if the outlier query is already well isolated. All the attributes identified by IPS-BS in SYN2 are correct (Prec=1), but only a small portion of true attributes are discovered (low Recall).
- The Recall scores are low for real-world datasets because we treat all original attributes to be outlying as ground truth. Here a low Recall does not necessarily indicate bad performance.

	COIN			CAL			IPS-BS			LIME		
	Prec	Recall	F1	Prec	Recall	F1	Prec	Recall	F1	Prec	Recall	F1
SYN1	0.97	0.89	0.93	0.89	0.81	0.84	0.87	0.44	0.58	0.82	0.79	0.80
SYN2	0.99	0.90	0.94	0.92	0.70	0.80	1.00	0.37	0.54	0.91	0.70	0.79
WBC	0.86	0.37	0.52	0.84	0.37	0.51	0.90	0.15	0.26	0.35	0.39	0.37
Twitter	0.91	0.33	0.48	0.75	0.34	0.47	0.72	0.29	0.41	0.60	0.67	0.63

Table 3: Faithfulness of Abnormal Attributes Identification

5.4 Outlierness Score Evaluation

Different outlier detection methods analyze data from different perspectives. Such differences will ultimately be reflected in the outlierness scores assigned to instances. A more effective detection mechanism is less likely to miss instances that are divergent from normal patterns, or consider normal instances to be more suspicious than true outliers. In this regard, the interpretation approach should be able to accurately measure the degree of deviation of a test outlier from its normal counterpart. In order to simulate ground truth outlierness, for each dataset applied in this experiment, we randomly sample the same number of inliers as the outliers, and use both of them as queries fed into interpreters. The ground truth score is 1 each true outlier, and 0 for inlier samples. For each query instance, interpreters are asked to estimate its outlierness score. After that, we rank the instances in descending order with respect to their scores. True outliers are more isolated than normal instances by their nature. A trustworthy interpreter should be able to maintain the relative magnitude of scores among all instances, so true outliers should be assigned with higher scores than inliers.

We report the results in Table 4 with AUC as the evaluation metric. We did not get valid result from IPS_BS for MNIST dataset as it fails in dealing with data of high dimensions, so its performance is not applicable here. The proposed method is advantageous over the baseline methods. In general, LIME slightly outperforms CAL. For SYN1 and WBC, the advantage of the proposed method is less obvious than that for other datasets. This can be explained by the differences of structural complexity among different datasets. For SYN1, an outlier is only detached from only one major cluster. For WBC, a malignant instance is usually characterized by those attributes with values significantly larger than normal. The contexts for these two datasets are relatively clear. However, for SYN2, Twitter and MNIST datasets, an outlier may be close to several separated neighboring clusters, thus producing outlier and inlier classes that are not trivially linear separable. Therefore, COIN and IPS_BS, once applicable, are more effective in these cases. IPS_BS is robust to complicated data structures, though it is less efficient than other methods. COIN resolves the context of outliers into clusters, so it can handle data of complex structures. It is worth noting that, LIME is more sensitive to model parameters, as it requires a complex global model upon which a set of local models are superimposed.

5.5 Interactions between Outlying Attributes and Outlierness

In real-world scenarios, outlier detection may serve for some practical purposes, such as spammer detection, fraud detection and

AUC	SYN1	SYN2	WBC	Twitter	MNIST
COIN	0.78	0.93	0.96	0.85	0.87
CAL	0.71	0.63	0.94	0.81	0.76
IPS_BS	0.69	0.90	0.90	0.79	0.74
LIME	0.74	0.62	0.94	0.83	0.78

Table 4: Outlierness score ranking performance

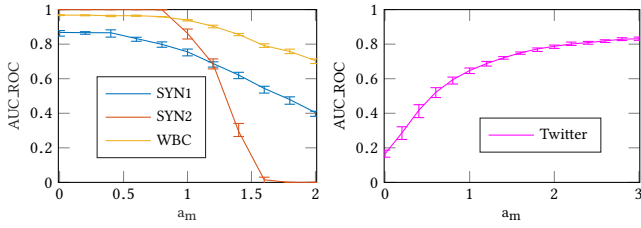
health monitoring. From the outlying attributes revealed by interpretation models, based on human knowledge, we can judge if their roles or semantics are in accordance with the nature of the problem. For those outliers whose abnormal attributes are loosely related to the problem, we want to weaken their significance or even discard them. In this experiment, we discuss how to refine the outlier detection results in terms of increasing the relevancy between spotted outliers and applications, by incorporating prior knowledge of the practical meaning of attributes.

The experiment is separated into two parts. In the first part, we assume that all the original attributes are equally relevant to the problem of interest, while some *simulated attributes* are appended to each instance. These attributes may cause new outliers to appear, but they are irrelevant to the ground truth. Similar to the previous experiment, we randomly sample the same number of inliers as test instances in addition to the true outliers. Here we set the number of simulated attributes to be the same as original ones, so each instance is augmented as $\mathbf{x} \in \mathbb{R}^{2M}$. We run COIN on different significance vectors β and set all entries in \mathbf{p} to be zero. The weights corresponding to original attributes are fixed to 1 ($\beta_m = 1, m \in [1, M]$), and we only vary the weights of simulated attributes ($\beta_m = \beta, m \in [M+1, 2M]$). Similar to Section 5.4, we obtain the outlierness score for all queries and rank them in descending order according to the score magnitude. True outliers are expected to have higher ranks than inliers. The performance of outlierness ranking is reported in Figure 4a. The plot indicates that as we increase the weights of simulated attributes, the performance of the interpreter degrades to varying degrees for all datasets, because it is more difficult for the interpreter to distinguish between real outliers and noisy instances. The degradation is not dramatic even when original and simulated attributes are weighted equally ($a_m = 1, m \in [M+1, 2M]$), which indicates that COIN is relatively robust to noisy data. However, as we increasingly misplace trust on simulated attributes that are irrelevant to the true outliers, factitious outlying instances start to dominate.

The second part of the experiment uses Twitter dataset which consists of the information of a number of normal users and spammers. The features extracted from user profiles, posts and graph

Selection	SYN1		SYN2		WBC		Twitter	
	d_{avg}	p_{real}	d_{avg}	p_{real}	d_{avg}	p_{real}	d_{avg}	p_{real}
COIN	0.75	0.875	0.94	0.97	0.95	0.90	0.70	0.81
CAL	0.58	0.75	0.57	0.80	0.88	0.81	0.73	0.48
IPS-BS	0.58	-	0.80	-	0.89	-	0.72	-
LIME	0.54	0.67	0.63	0.86	0.84	0.73	0.41	0.71

Table 5: Accuracy of model selection using outlier interpretation methods.



(a) Data with noise attributes (b) Twitter spammer data

Figure 4: The influence of the prior knowledge on outlier score. Results averaged over 20 runs, bars depict 25-75%.

structures are used as attributes. According to [53], the robustness level varies for different attributes. Some attributes, such as the number of followers, hashtag ratio and reply ratio, can be easily controlled by spammers to make themselves look normal, so that they are of low robustness. Other attributes such as account age, API ratio and URL ratio are beyond their easy control due to the huge potential expense or human labor, so they have high robustness. In this experiment, we fix the weight of low-robustness attributes to 1, and vary the weight β_m of high-robustness attributes. The entries of \mathbf{p} are decided according to [53]. The remaining procedures are the same as the first part of experiment discussed above. The number of normal instance queries is the half of the real outliers. The result of outlierness ranking is reported in Figure 4b. The rising curve shows that as more emphasis is put on high-robust attributes, we are able to refine the performance of identifying spammers. The experiment result indicates that by resorting to the interpretation of detected outliers, we can gain more insights on their characteristics, and more accurately select those that are in accordance with the purpose of the application.

5.6 Model Evaluation from Interpretation

In this experiment, we demonstrate that the interpretations can be used for model selection from detector evaluation without relying on the ground truth labels. Sometimes, end users may need to know the performance of competing methods on an outlier detection problem, in order to choose the most effective one to deploy in real applications. In this module, we add some noise attributes to the instances. Noisy attributes are seen as irrelevant to ground truth outliers, so their significance weights are set to zero in COIN. Similar to the experiment in subsection 5.5, outlierness incurred by noise attributes are undesired. The outlier detectors applied here include LOF [8], One-Class SVM [42] and Isolation Forest [29]. The three approaches are of different types, and involve disparate

definitions of outliers and algorithms to get solution. Several detectors can be built from the same approach from different parameter settings. For each dataset with $|\mathcal{O}|$ ground truth outliers, we let detectors return $1.5 \times |\mathcal{O}|$ outlier candidates. On one hand, we evaluate the performance of detectors using AUC with ground truth. On the other hand, the candidates are fed into interpreters to get deeper insight. Interpreters return the outlying attributes and outlierness scores as explanations, which provide two perspectives for evaluating the performance of detectors. First, as the noise attributes are irrelevant to the ground truth, original attributes are expected to be the outlying attributes for real outliers. We use $p_{real} = \sum_{m \in [1, M]} |s(a_m)| / \sum_{m \in [1, 2M]} |s(a_m)|$, i.e. the ratio of absolute weights of real attributes, to represent their relative significance. Second, we use the average distance d_{avg} as another metric to represent the effectiveness of the detection result. Given two detectors for comparison, the one which gets higher p_{real} or d_{avg} from an interpreter will be regarded as better and will be selected.

We generate 18 outlier detectors for each dataset, and pair up every two detectors with at least a gap of 0.05 in AUC. Meanwhile, for each pair of detectors, the interpreter also provides p_{real} and d_{avg} as two comparisons. If an interpreter could correctly evaluate the performance of detectors, then the detector with a higher AUC score tends to have greater p_{real} and d_{avg} than its competitor in the pair. Therefore, we select the detector with higher p_{real} and d_{avg} respectively, and check if it is consistent with the detector of higher AUC score. The accuracy of picking the correct detector is shown in Table 5. COIN is consistently better than the baseline methods. The p_{real} values of IPS-BS are usually equal for all pairs of detectors, so its accuracy is not applicable here. The capability of model selection of p_{real} and d_{avg} varies according to the structural complexity of datasets. For example, in SYN2 where data points aggregate in multiple clusters, attribute selection tends to make better choices than the distance measure. The results demonstrate that the assessment provided by interpretation can indicate the detection quality to varying degrees.

5.7 Case Studies

At last, we conduct some case studies to intuitively present the outcome of different components in COIN. MNIST dataset is used here as images are easier to understand perceptually. The attributes fed into the interpreter are hidden features extracted by the RBM. The latent features learned from RBM can be seen as simple primitives that compose more complicated visual patterns. It is more suitable for interpretation than using raw pixels as it is in accordance with the cognitive habits of people, that we tend to use richer representations for explanation and action [25].

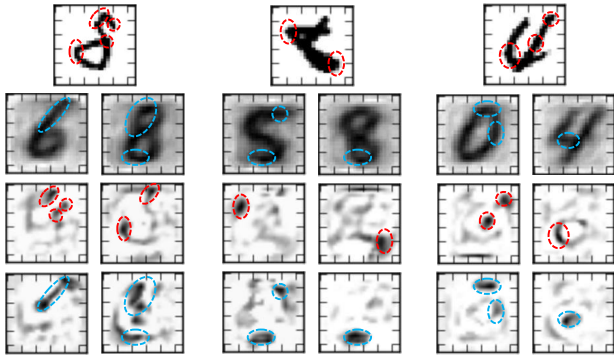


Figure 5: Examples of outlier interpretation using MNIST dataset. Red and blue circles highlight the regions explaining why images in the first row are recognized as outliers.

The case study results are shown in Figure 5. There are three query outlier images. The query outlier images are in the first row. We choose two neighboring clusters for each query, and obtain the average image of each cluster, as shown in the second row. Clear handwritten digits can be seen from average images, so that the clusters are internally coherent. The third and fourth rows indicate the characteristic attributes of the query image and average-neighbor image, respectively. The black strokes in the images of the third row represent positive outlying attributes, i.e., the query image is regarded as an outlier instance because it *possesses* these attributes. The strokes in fourth-row images are negative outlying attributes, as the query outlier digit does not include them. These negative attributes are, however, commonly seen in the neighbor images of certain cluster. The positive and negative attributes together explain why the outlier image is different from its nearby images in the dataset.

6 RELATED WORK

Many outlier detection approaches have been developed over the past decades. These approaches can be divided into three categories: density-based, distance-based and model-based approaches. Some notable density-based detection methods include [2, 8, 17, 44, 46]. Representative distance-based approaches include [4, 7, 22, 29, 40]. For model-based approaches, some well-known examples are [19, 42, 48]. Various approaches have been proposed to tackle the challenges including the curse of dimensionality [2, 16, 24], the massive data volume [4, 40], and heterogeneous information sources [17, 36]. Ensemble learning, which is widely used in supervised learning settings, can also be applied for outlier detection with non-trivial improvements in performance [28, 57]. [27] combines results from multiple outlier detectors, each of which apply only a subset of features. In contrast, each individual detector can subsample data instances to form an ensemble of detectors [57]. Some recent work starts to realize the importance about the explanations of detection results. In heterogeneous network anomaly detection, [17, 28, 30, 37] utilize attributes of nodes as auxiliary information for explaining the abnormality of resultant anomaly nodes. The motivation of this work is different from them, as we try to infer the reasons that why the given outliers are regarded as outlying, instead of developing new detection methods.

Besides algorithm development, researchers are also trying to provide explanations along with the approaches and their outcomes. The approach introduced in [21] can also find the subspace in which the features of outliers are exceptional. Ertöz *et al.* designed a framework for detecting network intrusion with explanations, which only works on categorical attributes [15]. The Bayesian program learning framework has been proposed for learning visual concepts that generalizes in a way similar to human, especially with just one or a few data examples [25]. Interpretations for anomalies detection can be naturally achieved within the scenario of attributed networks [17, 30, 37]. These techniques cannot be directly applied to solve our problem, because: (1) Heterogeneous information may not be available; (2) In many cases, features are not designed for achieving specific tasks; (3) The definition of anomalies varies in the work above, so a more general interpretation approach is still needed. Moreover, given the black-box characteristics of major mathematical models, the community is exploring ways to interpret the mechanisms that support the model, as well as the rules according to which the predictions are made. Ribeiro *et al.* developed a model-agnostic framework that infers explanations by approximating local input-output behavior of the original supervised learning model [41]. Lakkaraju *et al.* formalizes decision set learning which can generate short, succinct and non-overlapping rules for classification tasks [26]. Micenková *et al.* proposed to use classification models and feature selection methods to provide interpretations to the outliers in the subspace [33]. Vinh *et al.* utilize the isolation property of outliers and apply isolation forest for outlying aspects discovery [50].

7 CONCLUSION AND FUTURE WORK

In this paper, we propose the Contextual Outlier Interpretation (COIN) framework. The framework is model-agnostic and can be applied to a wide range of detection methods. The goal of interpretation is achieved by solving a series of classification tasks. Each outlier query is explained within its local context. The abnormal attributes and outlierness score of an outlier can be obtained by a collection of simple but interpretable classifiers built in its resolved context. We also propose a new measure of outlierness score whose relationship with abnormal attributes can be explicitly formulated. Prior knowledge on the roles of attributes in different scenarios can also be easily incorporated into the interpretation process. The explanatory information of multiple queries can be aggregated for evaluating detection models. Comprehensive evaluation on interpretation performance and model selection accuracy are provided through a series of experiments with both real world and simulated datasets. Case studies are also conducted for illustrating the outcome of each component of the framework.

There are a number of directions for future work that can be further explored. Hierarchical clustering strategies can be designed to more accurately resolve the context of an outlier query for better interpretation. The framework can be extended to handle heterogeneous data sources. Moreover, strategies for dealing with outlier groups can be designed, so that interpretation approaches can be applied to a wider range of objects.

REFERENCES

- [1] Naoki Abe, Bianca Zadrozny, and John Langford. 2006. Outlier detection by active learning. In *KDD*. ACM, 504–509.
- [2] Charu C Aggarwal and Philip S Yu. 2001. Outlier detection for high dimensional data. In *ACM Sigmod Record*, Vol. 30. ACM, 37–46.
- [3] Leman Akoglu, Mary McGlohan, and Christos Faloutsos. 2010. OddBall: Spotting anomalies in weighted graphs. *PAKDD* (2010).
- [4] Fabrizio Angiulli and Fabio Fasseti. 2009. Dolphin: An efficient algorithm for mining distance-based outliers in very large datasets. *TKDD* 3, 1 (2009), 4.
- [5] A. Asuncion and D.J. Newman. 2007. UCI Machine Learning Repository. (2007). <http://www.ics.uci.edu/~sim5mllearn/>
- [6] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert MÄzler. 2010. How to explain individual classification decisions. *JMLR* 11, Jun (2010), 1803–1831.
- [7] Stephen D Bay and Mark Schwabacher. 2003. Mining distance-based outliers in near linear time with randomization and a simple pruning rule. In *KDD*. ACM, 29–38.
- [8] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. LOF: identifying density-based local outliers. In *ACM sigmod record*, Vol. 29. ACM, 93–104.
- [9] Guilherme O Campos, Arthur Zimek, Jörg Sander, Ricardo JGB Campello, Barbora Micenkova, Erich Schubert, Ira Assent, and Michael E Houle. 2015. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data Mining and Knowledge Discovery* (2015), 1–37.
- [10] Davide Canali, Marco Cova, Giovanni Vigna, and Christopher Kruegel. 2011. Prophiler: a fast filter for the large-scale detection of malicious web pages. In *Proceedings of the 20th international conference on World Wide Web*.
- [11] Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41, 3 (2009), 15.
- [12] Jesse Davis and Mark Goadrich. 2006. The relationship between Precision-Recall and ROC curves. In *Proceedings of the 23rd international conference on Machine learning*. ACM, 233–240.
- [13] Lei Duan, Guanting Tang, Jian Pei, James Bailey, Guozhu Dong, Akiko Campbell, and Changjie Tang. 2014. Mining contrast subspaces. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 249–260.
- [14] Bradley Efron, Trevor Hastie, Iain Johnstone, Robert Tibshirani, and others. 2004. Least angle regression. *The Annals of statistics* 32, 2 (2004), 407–499.
- [15] Levent Ertöz, Eric Eilertson, Aleksandar Lazarevic, Pang-Ning Tan, Vipin Kumar, Jaideep Srivastava, and Paul Dokas. 2004. Minds-minnesota intrusion detection system. *Next generation data mining* (2004), 199–218.
- [16] Peter Filzmoser, Ricardo Maronna, and Mark Werner. 2008. Outlier identification in high dimensions. *Computational Statistics & Data Analysis* 52, 3 (2008).
- [17] J. Gao, F. Liang, W. Fan, C. Wang, Y. Sun, and J. Han. 2010. On community outliers and their efficient detection in information networks. *KDD* (2010).
- [18] Haibo He and Edwardo A Garcia. 2009. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering* 21, 9 (2009), 1263–1284.
- [19] Zengyou He, Xiaofei Xu, and Shengchun Deng. 2003. Discovering cluster-based local outliers. *Pattern Recognition Letters* 24, 9 (2003), 1641–1650.
- [20] Geoffrey E Hinton and Ruslan R Salakhutdinov. 2006. Reducing the dimensionality of data with neural networks. *science* 313, 5786 (2006), 504–507.
- [21] Edwin M Knorr and Raymond T Ng. 1999. Finding intensional knowledge of distance-based outliers. In *VLDB*, Vol. 99. 211–222.
- [22] Edwin M Knorr, Raymond T Ng, and Vladimir Tucakov. 2000. Distance-based outliers: algorithms and applications. *The VLDB Journal/The International Journal on Very Large Data Bases* 8, 3-4 (2000), 237–253.
- [23] Hans-Peter Kriegel, Peer Kröger, Erich Schubert, and Arthur Zimek. 2011. Interpreting and Unifying Outlier Scores.. In *SDM*. SIAM, 13–24.
- [24] Hans-Peter Kriegel, Peer Kröger, and Arthur Zimek. 2009. Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 3, 1 (2009), 1.
- [25] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. 2015. Human-level concept learning through probabilistic program induction. *Science* 350, 6266 (2015), 1332–1338.
- [26] H Lakkaraju, S Bach, and J Leskovec. 2016. Interpretable decision sets: A joint framework for description and prediction. In *KDD*. ACM, 1675–1684.
- [27] Aleksandar Lazarevic and Vipin Kumar. 2005. Feature bagging for outlier detection. In *KDD*. ACM, 157–166.
- [28] Jiongqian Liang and Srinivasan Parthasarathy. 2016. Robust contextual outlier detection: Where context meets sparsity. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*.
- [29] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2012. Isolation-based anomaly detection. *TKDD* 6, 1 (2012), 3.
- [30] Ninghao Liu, Xiao Huang, and Xia Hu. 2017. Accelerated Local Anomaly Detection via Resolving Attributed Networks. *IJCAI*.
- [31] Yuli Liu, Yiqun Liu, Ke Zhou, Min Zhang, and Shaoping Ma. 2017. Detecting Collusive Spamming Activities in Community Question Answering. In *Proceedings of the 26th International Conference on World Wide Web*.
- [32] Junshui Ma and Simon Perkins. 2003. Online novelty detection on temporal sequences. In *KDD*. ACM, 613–618.
- [33] Barbora Micenkova, Raymond T Ng, Xuan-Hong Dang, and Ira Assent. 2013. Explaining outliers by subspace separability. In *2013 IEEE 13th International Conference on Data Mining*. IEEE, 518–527.
- [34] Emmanuel Muller, Patricia Iglesias Sánchez, Yvonne Mulle, and Klemens Bohm. 2013. Ranking outlier nodes in subspaces of attributed graphs. *ICDEW* (2013).
- [35] Alexandros Ntoulas, Marc Najork, Mark Manasse, and Dennis Fetterly. 2006. Detecting spam web pages through content analysis. In *Proceedings of the 15th international conference on World Wide Web*.
- [36] Bryan Perozzi and Leman Akoglu. 2016. Scalable anomaly ranking of attributed neighborhoods. *SDM* (2016).
- [37] Bryan Perozzi, Leman Akoglu, Patricia Iglesias Sánchez, and Emmanuel Müller. 2014. Focused clustering and outlier detection in large attributed graphs. *KDD* (2014).
- [38] Clifton Phua, Vincent Lee, Kate Smith, and Ross Gayler. 2010. A comprehensive survey of data mining-based fraud detection research. *arXiv preprint arXiv:1009.6119* (2010).
- [39] Leonid Portnoy, Eleazar Eskin, and Sal Stolfo. 2001. Intrusion detection with unlabeled data using clustering. In *In Proceedings of ACM CSS Workshop on Data Mining Applied to Security (DMSA-2001)*. Citeseer.
- [40] S. Ramaswamy, R. Rastogi, and K. Shim. 2000. Efficient algorithms for mining outliers from large data sets. In *ACM SIGMOD Record*, Vol. 29. ACM, 427–438.
- [41] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *KDD*.
- [42] Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. 2001. Estimating the support of a high-dimensional distribution. *Neural computation* 13, 7 (2001), 1443–1471.
- [43] Neil Shah. 2017. FLOCK: Combating Astroturfing on Livestreaming Platforms. In *Proceedings of the 26th International Conference on World Wide Web*.
- [44] Xiuyao Song, Mingxi Wu, Christopher Jermaine, and Sanjay Ranka. 2007. Conditional anomaly detection. *IEEE Transactions on Knowledge and Data Engineering* (2007).
- [45] J. Sun, H. Qu, D. Chakrabarti, and C. Faloutsos. 2005. Neighborhood formation and anomaly detection in bipartite graphs. *ICDM* (2005).
- [46] Jian Tang, Zhixiang Chen, Ada Wai-Chee Fu, and David W Cheung. 2002. Enhancing effectiveness of outlier detections for low density patterns. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 535–548.
- [47] Robert Tibshirani and Guenther Walther. 2005. Cluster validation by prediction strength. *Journal of Computational and Graphical Statistics* 14, 3 (2005), 511–528.
- [48] Hanghang Tong and Ching-Yung Lin. 2011. Non-Negative Residual Matrix Factorization with Application to Graph Anomaly Detection. *SDM* (2011).
- [49] Simon Tong and Daphne Koller. 2000. Restricted bayes optimal classifiers. In *AAAI/IAAI*. 658–664.
- [50] Nguyen Xuan Vinh, Jeffrey Chan, Simone Romano, James Bailey, Christopher Leckie, Kotagiri Ramamohanarao, and Jian Pei. 2016. Discovering outlying aspects in large datasets. *Data Mining and Knowledge Discovery* (2016), 1–36.
- [51] W. Wong, A. Moore, G. Cooper, and M. Wagner. 2002. Rule-based anomaly pattern detection for detecting disease outbreaks. In *AAAI/IAAI*. 217–223.
- [52] Liang Xiong, Xi Chen, and Jeff Schneider. 2011. Direct robust matrix factorization for anomaly detection. *ICDM* (2011).
- [53] Chao Yang, Robert Chandler Harkreader, and Guofei Gu. 2011. Die free or live hard? empirical evaluation and new design for fighting evolving twitter spammers. In *International Workshop on Recent Advances in Intrusion Detection*. Springer, 318–337.
- [54] Yiming Yang, Jian Zhang, Jaime Carbonell, and Chun Jin. 2002. Topic-conditioned novelty detection. In *KDD*. ACM, 688–693.
- [55] Zhe Zhao, Paul Resnick, and Qiaozhu Mei. 2015. Enquiring minds: Early detection of rumors in social media from enquiry posts. In *Proceedings of the 24th International Conference on World Wide Web*.
- [56] Ji Zhu, Saharon Rosset, Trevor Hastie, and Rob Tibshirani. 2004. 1-norm support vector machines. *NIPS* 16, 1 (2004), 49–56.
- [57] Arthur Zimek, Matthew Gaudet, Ricardo JGB Campello, and Jörg Sander. 2013. Subsampling for efficient and effective unsupervised outlier detection ensembles. In *KDD*. ACM, 428–436.