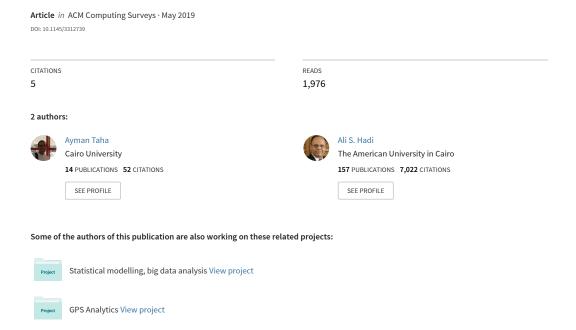
See discussions, stats, and author profiles for this publication at: https://www.researchgate.net/publication/333524837

Anomaly Detection Methods for Categorical Data: A Review



Anomaly Detection Methods for Categorical Data: A Review

AYMAN TAHA, Faculty of computers and Information- Cairo University, Egypt. a.taha@fci-cu.edu.eg ALI S. HADI, American University in Cairo, Egypt, and Cornell University, Ithaca, NY, USA. ahadi@aucegypt.edu or ali-hadi@cornell.edu

Anomaly detection has numerous applications in diverse fields. For example, it has been widely used for discovering network intrusions and malicious events. It has also been used in numerous other applications such as identifying medical malpractice or credit fraud. Detection of anomalies in quantitative data has received a considerable attention in the literature and has a venerable history. By contrast, and despite the widespread availability use of categorical data in practice, anomaly detection in categorical data has received relatively little attention as compared to quantitative data. This is because detection of anomalies in categorical data is a challenging problem. Some anomaly detection techniques depend on identifying a representative pattern then measuring distances between objects and this pattern. Objects that are far from this pattern are declared as anomalies. However, identifying patterns and measuring distances are not easy in categorical data compared with quantitative data. Fortunately, several papers focussing on the detection of anomalies in categorical data have been published in the recent literature. In this article, we provide a comprehensive review of the research on the anomaly detection problem in categorical data. Previous review articles focus on either the statistics literature or the machine learning and computer science literature. This review article combines both literatures. We review 36 methods for the detection of anomalies in categorical data in both literatures and classify them into 12 different categories based on the conceptual definition of anomalies they use. For each approach, we survey anomaly detection methods, and then show the similarities and differences among them. We emphasize two important issues, the number of parameters each method requires and its time complexity. The first issue is critical because the performance of these methods are sensitive to the choice of these parameters. The time complexity is also very important in real applications especially in big data applications. We report the time complexity if it is reported by the authors of the methods. If not we derive it ourselves and report it in this paper. In addition, we discuss the common problems and the future directions of the anomaly detection in categorical data.

CCS Concepts: • Information systems applications \rightarrow Data mining.

Additional Key Words and Phrases: Computational complexity, Data mining, Holo entropy, Intrusion detection systems, Mixed data, Outliers detection, Nominal data, Novelty detection, Semi-supervised learning, Shannon entropy, Supervised learning, Unsupervised learning

ACM Reference Format:

Ayman Taha and Ali S. Hadi. 2019. Anomaly Detection Methods for Categorical Data: A Review. 1, 1 (January 2019), 35 pages. https://doi.org/10.1145/1122445.1122456

Authors' addresses: Ayman Taha, Faculty of computers and Information-Cairo University, Egypt. a.taha@fci-cu.edu.eg; Ali S. Hadi, American University in Cairo, Egypt, and Cornell University, Ithaca, NY, USA. ahadi@aucegypt.edu or alihadi@cornell.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

XXXX-XXXX/2019/1-ART \$15.00

https://doi.org/10.1145/1122445.1122456

1 INTRODUCTION

Anomalies are a minority of objects (observations, cases, or points) that are inconsistent with the pattern suggested by the majority of objects in the same dataset. It is very important to identify anomalies in data because they can distort the analysis and the decisions based on the analysis [23], [43], [81], [82], [85], [83], and [30]. Anomaly detection (or outlier or novelty detection as it is sometimes referred to in the literature) is a well studied problem in different research areas such as communications, statistics, data mining and machine learning (see, e.g., [80], [33], [16], [151], [97], [115], and [45]).

The massive increases in e-commerce applications, cloud computing services and remote business access have triggered the need for protection systems against unauthorized access. In communications, illegal access to network resources is called intrusion. Intrusion Detection Systems (IDS) are substantial parts in network security that observe computer networks to detect odd activities [53]. IDS are classified into two main categories: Signature-based and Anomaly-based.

A signature-based intrusion detection system relies on building a database of defined signatures for known attacks. It raises an alert when observing signature similar to those attacks in its database. It fails to detect new attacks whose signatures do not exist in the database. On the other hand, an anomaly-based intrusion detection system usually builds a statistical model defining normal activities. Abnormal activities (anomalies) are those activities that are discordant with the defined model [181].

Real datasets often consist of different types of variables, that is, some variables are quantitative and others are qualitative or categorical. These data are called mixed data. The analysis of quantitative data has a venerable history and, by comparison, categorical data has received less attention than quantitative data. Categorical data, however, are common in many different domains (e.g., biomedical, educational, sociological, psychological, political, and social sciences ([7], [8]). Anomaly detection techniques for categorical data make use of data within these fields [172] and [15].

There are two types of categorical data: nominal and ordinal. Examples of the former type include gender, nationality, and type of network protocol. Examples of the latter type include letter grade in a course (e.g., A, B, C, D, F), network traffic volume (e.g., low, medium, high), and a Likert scale variable (e.g., 1 = strongly disagree, 2 = disagree, 3 = neutral, 4 = agree, and 5 = strongly agree). Finding anomalies in ordinal variables is relatively easier than the identification of anomalies in nominal data because the categories in nominal variables have no natural ordering [174].

1.1 Applications of Anomaly Detection in Categorical Data

Applications of anomaly detection in categorical data are numerous. We mention here just few examples:

- Network Intrusions: Detecting users who have abnormal accesses is an interesting application for anomaly detection in cloud computing [154]. Here users' profiles and behaviors are usually coded as categorical variables.
- Moving Objects: Anomaly detection has useful applications in analyzing moving objects data such as identifying objects that have unexpected moves and identifying types of road segments, where moving objects have anomalous behavior [71] and [116].
- Medical and Health Data: Anomaly detection methods provide useful information to decision makers such as the identification of excessively high medical payments and/or medical malpractice. This information can be used to improve health care management [96]).

- Social Networks: There are many applications that make use of anomaly detection in social networks such as finding anomalous users in social groups who have different viewpoints, interests, and believes [5], and [69].
- Credit Fraud: Unauthorized credit card use may be identified when an unexpected behavior is detected. Categorical attributes are often used to represent behavioral data [175].
- Questionnaire Data: Responses to questionnaires are usually mapped into categorical or ordinal attributes. The importance of anomaly detection in questionnaire data is highlighted in [187].
- Earth Science: Finding anomalies in spatiotemporal data, e.g., weather patterns or climate changes in various geographical areas give an explanation behind interesting spatiotemporal patterns [117].
- Law Enforcement: Examples of anomaly detection applications for law enforcement are discovering anomalies in trading activities and insurance claims [3].

1.2 Challenges Facing Anomaly Detection in Categorical Data

The identification of anomalies in categorical data faces some challenges. These include:

- (1) Anomaly detection methods focus on identifying the pattern (statistical distribution) suggested by the majority of observations, then considering the observations that do not follow the assumed pattern as anomalies [30]. Some distance functions are proposed in the literature to compute distance between categorical observations (e.g., [32] and [41]). Identifying a pattern and measuring a distance are not easy in categorical data. Consequently, anomaly detection methods are more common for quantitative than for categorical data.
- (2) Several alternative but different definitions of anomalies in categorical data exist in the literature [179]. Anomaly detection methods can identify different sets of observations as anomalies depending on the definition they adopt.
- (3) The number of benchmark datasets that can be used to test the performance of anomaly detection methods for categorical data (e.g. computation time, detection rate, etc.) is very small. In addition, it is not straightforward to generate synthetic categorical data with known anomalies because of the lack of methods that generate such data [171].
- (4) Computational complexity is also a challenging problem in the identification of anomalies especially in categorical data because most of real applications have huge datasets in terms of the number of observations, the number of categorical variables, and the number of categories in each. Consequently, time complexity is a significant issue when applying anomaly detection methods to categorical data.

1.3 Related Work

The problem of anomaly detection has been reviewed in a number of survey articles, books and book chapters (see, e.g., [122], [123], [91], [38], [39], [100], [80], and [46]). A recent book in [3], one can find an account of various outliers detection approaches for quantitative as well as categorical data. Chapter 8 of this book in particular discusses the problem of identifying outliers in categorical data.

Some survey articles are generic, giving a broad overview of anomaly detection methods (see, e.g., [163], [72], [59], [119], and [60]). Other review articles focus only on one type of methods based on application, technical approach and/or data type. These methods can be classified according to application domain as follows:

• Network Intrusion Detection: These articles survey, compare, classify and/or identify challenges in intrusion detection methods in computer networks (see, e.g., [61], [108], [134],

- [70], [75], [78], [20], and [10]). Distance and similarity measures used in intrusion detection methods are surveyed in [177].
- Wireless Sensor Network (WSN): This group of articles focus on anomaly detection techniques for wireless sensor data (see, e.g., [185], [145], [49], [118], [149], [1], [146], [2], [161], and [121]). Characteristics of anomaly detection techniques for wireless sensor networks in a non-stationary and harsh environments are surveyed in [130] and [159]. Anomaly detection in smart home and automated surveillance are reviewed in [165] and [21].
- Data Streams: Few review articles study the features and classification of anomaly detection in data stream (see, e.g., [19], [66], and [55]). Anomaly detection methods for dynamic streaming data are reviewed in [184]. Also, [141] focus on anomalies in continuous time variant data stream.
- Fraud Detection: Fraud and abuse detection methods have useful applications in many domains. Health care fraud detection techniques are surveyed in [99], [48], and [99]. Fraud detection methods for biological and chemical data are reviewed in [46]. Fault detection methods in industrial processes are studied in [27].
- Financial, Business and Recommender Systems: Financial and credit card fraud detection methods are surveyed in [138], [135], [101], [11], and [178]. Outlying profile attacks in recommender systems are studied in [58].

Another set of survey articles focus on the approach or the methodologies that are used to identify outliers. These approaches include:

- Data Mining Approach: There are several review articles that survey data mining-based anomaly detection methods (see, e.g., [61], [26], [99], [102], and [12]). Other review articles cover a certain data mining pattern (such as neural network-based novelty detection methods are covered in [123]), while, frequent pattern-based anomaly detection methods are reviewed in [18]. Clustering-based anomaly detection methods are reviewed in [55]. Anomaly detection techniques for distributed data are studied in [12]. Frequent pattern-based anomaly detection methods and their scoring measures are discussed in [153].
- Machine Learning Approach: There are few survey articles that cover machine learning-based anomaly detection methods (see, e.g., [139], [27], and [20]).
- Statistical Approach: The anomaly detection problem appeared in the statistical literature before data mining and machine learning. Statistical anomaly detection methods concentrate on simple data types (such as numerical and quantitative data) rather than complex structures (such as categorical, graphical and/or spatial data). Statistical anomaly detection methods are surveyed in [122], [84], [72], [52], and [46].

A third group of review articles concentrate on the type of data where anomaly detection methods are proposed to identify outliers within this type of data. These articles can be classified as follows:

- Social Networks and Graph Data: Anomaly detection methods in social media are compared and reviewed in some survey articles (see, e.g., [148], [156], [102], and [183]). However, some surveys focus on wider scope, graph data (see, e.g., [155], [143], and [14]).
- Time Series and Spatial Data: Several anomaly detection methods are proposed for time series and spatial data (see, e.g., [19], [80], [79], and [141]).
- Complex and Big Data: Big and complex data are active areas of research particularly in the computer science literature. Recently, several methods are proposed for the identification of anomalies in high dimensional, complex, and big data (see, e.g., [57] and [65]). Anomaly detection for high-dimensional data are given in [4] and [188]. Other techniques focus on the identification of anomalies in subspaces rather than studying all dimensions [166].

• Empirical Comparative Reviews: Some articles focus on experimental comparative study among a certain type of anomaly detection methods such as some statistical and distance-based anomaly detection methods (see, e.g., [22] and [132]). Further, empirical evaluations of unsupervised anomaly detection measures are conducted in [36].

Most of the survey articles for anomaly detection provide extensive surveys for detection of anomalies in quantitative data (see, e.g., [42], [91], [9], [84], [38], and [100]). Some of these surveys mention categorical data but just in passing. Discrete sequence anomaly detection methods are reviewed in [39].

1.4 Contributions

This article provides a comprehensive review of the research on anomaly detection in categorical data. Particular contributions are:

- (1) The interest in the identification of anomalies in categorical data is recent, but fortunately it has been intensifying during the last decade. To the best of our knowledge, this is the first review paper that focuses on the identification of anomalies in categorical data.
- (2) The papers on the identification of anomalies in categorical data appear in both the statistical literature as well as the computer science literature such as machine learning and big data. This paper reviews both of these literatures.
- (3) The existing anomaly detection methods for categorical data are categorized into different approaches. Each anomaly detection method is reviewed and discussed. Discussions of the number of parameters they require and comparisons among these methods are provided in Sections 2 to 10.
- (4) Since computational complexity is a significant issue in real applications, for each method, we report the time complexity if it is reported by the author(s) of the methods. If not, we derive and report it in this article. Tables 1 and 2 in the Appendix show the reviewed anomaly detection methods together with their computational complexities and required parameters.
- (5) Extensive discussion of the problems facing anomaly detection methods in categorical data are provided and topics for future research in anomaly detection in categorical data are presented.

1.5 Organization

The rest of this paper is organized as follows: Sections 2–10 give a review of anomaly detection methods in categorical data. Discussions are provided in Section 11. Section 12 provides concluding remarks and recommendations for future work. Tables 1 and 2 in the Appendix summarize the reviewed anomaly detection methods for categorical data, their complexities, and their required parameters.

2 METHODS BASED ON INDICATOR VARIABLES

One approach for the detection of anomalies in categorical data is to represent categorical data by numeric values. Then anomaly detection for quantitative data can be used. Examples of methods that follow this approach are (a) Indicator Variables (IV) [162] and (b) A method based on Multiple Correspondence Analysis (MCA) [162].

The IV method replaces each categorical variable by indicator variables [162]. For example, a categorical variable X_j with c_j categories, can be represented by c_j indicator (binary) variables, $Y_{n \times c_j} = \{Y_1^j, Y_2^j, \cdots, Y_{c_j}^j\}$, where Y_r^j is the r-th indicator variable corresponding to the r-th category of the j-th categorical variable X_j . To identify the anomalies in this binary dataset with Q indicator variables, where Q is the total number of categories in the dataset, the Canberra distance [64] is

used. For each observation x_i in the binary dataset, it computes an outlying score, S_i , as the average of the Canberra distance between x_i and all other n-1 observations in the dataset. It declares an observation x_i as an anomaly if $S_i > \theta$, where θ is a predefined parameter chosen by the user [162].

The MCA is a generalization of Correspondence Analysis (CA) when the number of categorical variables is greater than two, [25]. The MCA-based anomaly detection method represents categorical variables by indicator variables [162] and use the Canberra distance to identify the anomalies in the transformed numerical datasets. The MCA steps are then applied to the binary matrix B of size $n \times \sum_{j=1}^{q} c_j$ instead of the contingency table, which is used by the CA.

The time complexity of the Indicator Variables method can be divided into two parts. The time complexity of creating the indicator variables, which is O(nq), and the time complexity of computing the Canberra distance, which is $O(n^2 \sum_{j=1}^q (c_j - 1))$. Accordingly, the total time complexity for the Indicator Variable method is $O(nq + n^2 \sum_{j=1}^q (c_j - 1))$.

The time complexity of the MCA-based method consists of three parts. First, the time complexity of creating the indicator variables is O(nq). Second, the time complexity of finding the singular value decomposition is $(nQ^2 + n^2Q + Q^3)$ ([76] and [54]). Finally, the time complexity of computing the Canberra distance is $O(n^2Q)$. Therefore, the total time complexity is $O(nq + 2n^2Q + nQ^2 + Q^3)$.

As can be seen, the Indicator Variable method is computationally expensive due to the substantial increase of number of indicator variables and the quadratic cost of measuring distance with respect to the number of observations. The use of indicator variables method is not recommended when the number of categories is large. The authors solve this problem by grouping similar categories into one category using some clustering methods. However, this solution slightly decreases the number of indicator variables but it leads to loss of information due to grouping [162].

The MCA-based anomaly detection method is also computationally expensive due to the computation of singular value decomposition in addition to the quadratic cost of measuring distance with respect to the number of observations.

3 FREQUENCY-BASED METHODS

Frequency-based methods make use of the frequency (i.e., the number of occurrences of categories) instead of distances. There are three types of frequencies that can be used to identify anomalies in categorical data: (a) Marginal frequency, (b) Itemset Frequency, and (c) Diversified Frequency.

3.1 Marginal Frequency

Anomalies in categorical data can be defined as observations with small marginal frequency. Methods using this definition are (a) the Attribute Value Frequency (AVF) [106] and [105], (b) the Square of the Complement Frequency (SCF) [173], (c) the Weighted Attribute Value Frequency (WAVF) [150], (d) the Weighted Density-based Outlier Detection (WDOD) method [186], (e) a Cloud Model-Based Outlier Detection method (CMBOD) [109], and (f) Bouguessa's method [34].

The AVF is based on computing a frequency score for each observation, x_i , as

$$AVF(x_i) = \frac{1}{q} \sum_{j=1}^{q} f(x_{ij}), \tag{1}$$

where $f(x_{ij})$ is the marginal frequency of x_{ij} within the variable X_j . The AVF declares the M objects with the lowest AVF scores as outliers, where M is a parameter chosen by the user. Alternatively, [147] suggest estimating the number of anomalies M by assuming that the frequency scores follow a normal distribution.

The Square of the Complement Frequency (SCF) calculates an outlying score for each observation, x_i , in the dataset as

$$SCF(x_i) = \sum_{j=1}^{q} \frac{[1 - p(x_{ij})]^2}{c_j},$$
(2)

where c_j is the number of categories in the variable X_j and $p(x_{ij})$ is the marginal relative frequency of x_{ij} (the number of occurrence within X_j divided by number of observations). It then declares the M observations with the highest outlying scores as outliers. However, SCF uses the square of the complement frequency to increase the difference between frequent and infrequent categories. Moreover, SCF takes into consideration the number of categories (c_j) to give variables with small c_j higher weights in the score function.

The AVF and SCF methods do not take into account the sparseness of the frequencies in categorical variables. Therefore the performance of these methods can be improved by capturing the sparsity of categorical variables by giving weight to the frequencies. Two examples of such weight are (a) the Weighted Attribute Value Frequency (WAVF) [150] and (b) the Weighted Density-based Outlier Detection (WDOD) method [186].

In the WAVF a weighting function is presented to indicate the variable sparsity in measuring the outlying scores. The higher sparsity the variable has, the higher impact on the outlying score. Thus observations with more sparse categories have higher probability for being outliers. The sparseness level of a categorical variable can be measured by sparsity statistical functions, e.g., the range or the standard deviation of the marginal frequencies. The Weighted Attribute Value Frequency (WAVF) is proposed to extend the AVF method as

WAVF
$$(x_i) = \frac{1}{q} \sum_{j=1}^{q} f(x_{ij}) * R_j,$$
 (3)

where R_j is the range of the frequencies of the jth categorical variable.

Another way of giving weights to the frequencies is the weighted density-based outlier detection (WDOD) method for categorical data [186]. It relies on estimating the density of each categorical variable then it computes the weighted density for the whole dataset. More specifically, WDOD starts with computing a weight, $W(X_j)$, for each variable X_j as

$$W(X_j) = \frac{1 - E^c(X_j)}{\sum_{l=1}^{q} (1 - E^c(X_l))},$$
(4)

where

$$E^{c}(X_{j}) = \sum_{i=1}^{c_{j}} p(x_{ij}) * (1 - p(x_{ij})),$$
 (5)

is the complement entropy for X_j , $p(x_{ij})$ is the relative frequency of a category x_{ij} , and c_j is number of categories in X_j . According to [186], the complement entropy, $E^c(X_j)$, is proposed in [114] to measure the information gain or uncertainty to give different weights to categorical variables to emphasize their importance. Unlike the logarithmic behavior of the Shannon entropy, the complement entropy measures both the uncertainty and the fuzziness. The WDOD score for an object x_i is then obtained by the sum of the weighted relative frequencies, that is,

$$WDOD(x_i) = \sum_{j=1}^{q} p(x_{ij}) * W(X_j).$$
(6)

Accordingly, WDOD(x_i) takes into consideration the uncertainty in the density of each variable. The lower weighted density WDOD(x_i) an object has, the higher probability of being a density-based outlier. Therefore, the WDOD declares an observation, x_i , as an outlier if (WDOD(x_i) < θ), where θ is a predefined parameter.

The Cloud Model-Based Outlier Detection method (CMBOD) [109] is based on building a cloud model which summarizes a dataset into a multidimensional cloud model [110]. The method consists of three steps. Step 1 transforms categorical variables to numerical values through a probabilistic mapping method called *cloud drops*. For a categorical variable X_j with c_j categories, each category is replaced by its marginal relative frequency $p(x_{ij})$.

Step 2 uses the transformed dataset in building multidimensional cloud model as follows: Three digital characteristics of the dataset are extracted. The digital characteristics of a variable X_j are (a) the average of X_j , \bar{x}_j , (b) the entropy of X_j , E_{nj} , and (c) the hyper entropy of X_j , E_{hj} , [111]. The entropy of an observation is calculated as $E_{nj}(x_{ij}) = \sqrt{\frac{\pi}{2}} * |(x_{ij} - \bar{x}_j)|$. The entropy of a variable X_j is the average entropies of all observations, which is

$$E_{nj} = \frac{1}{n} \sum_{i=1}^{n} E_{nj}(x_{ij}), \quad j = 1, 2, \cdots, q.$$
 (7)

The hyper entropy, E_{hj} , of a dataset is the average measure for the fussiness of entropy values. It is calculated as $E_{hj} = \sqrt{S_j^2 - E_{nj}^2}$ [86], where

$$S_j^2 = \frac{1}{n-1} \sum_{i=1}^n \left(x_{ij} - \bar{x}_j \right)^2, \quad j = 1, 2, \cdots, q,$$
 (8)

is the variance of X_i .

Step 3 computes the certainty degree for each observation x_i as:

$$cd(x_i) = \exp\left(\sum_{j=1}^{q} \frac{-(x_{ij} - \bar{x}_j)^2}{2E_{nj}^2}\right), \quad i = 1, 2, \dots, n.$$
 (9)

The certainty degree $cd(x_i)$ is used as a measure of belonging to the established cloud model. The lower certainty degree means higher deviation from the established cloud model, that is, the higher probability of being an anomaly. The CMBOD labels the M observations with the smallest $cd(x_i)$ as outliers.

The method proposed in [34] is designed to detect anomalies in mixed datasets. It is based on probabilistic approach that employs univariate beta mixture model (the bivariate beta mixture model in the case of mixed data) to distinguish outlying from non-outlying data. The method computes two outlying scores for each observation, one score for the quantitative attributes and the other score for the categorical attributes.

The outlying score of the p quantitative (numerical) variables is

$$S_N(x_i) = \sum_{j=1}^p \log \left(1 + \sum_{r=1}^k [x_{ij} - \text{knn}_r(x_{ij})]^2 \right), \tag{10}$$

where p is number of quantitative (numerical) variables and $\operatorname{knn}_r(x_{ij})$ is the value of the r-th nearest neighbor to (x_{ij}) in the j^{th} quantitative variable. Whereas the outlying score of the q categorical part is

$$S_C(x_i) = \sum_{j=1}^q \log \left(f(x_{ij}) \right), \tag{11}$$

where q is the number of categorical variables and $f(x_{ij})$ is the frequency of (x_{ij}) in the j^{th} categorical variable. The outlying scores use the log to emphasize the difference between frequent and rare values. Since large values of $S_N(x_i)$ indicates outliers but small values of $S_C(x_i)$ indicate outliers. Therefore, $S_C(x_i)$ is reversed and replaced by

$$S_C(x_i) \leftarrow \max(S_C(x_i)) - S_C(x_i) \tag{12}$$

so that large values of $S_C(x_i)$ now indicate outliers. Each of the two scores are then normalized to have values between 0 and 1. A bivariate beta distribution can then be used to determine if the normalized scores are large enough to be declared as outliers.

The author claims that this method does not require any decision parameters but this claim is not entirely true, since it requires three decision parameters. It is true that the author avoids requiring the number of outliers, M, by using a mixture of m univariate beta components. Then, it labels observations that corresponds to the vectors with the highest scores as outliers. It estimates number of components, m, by trying each possible values of m from $1, 2, \cdots, M_{\max}$, which is a user predefined parameter. In addition, in each iteration, it partitions the scores by the K-means method then it applies Expectation Maximization (EM) Algorithm to estimate the parameters of the m components. Thus, the second parameter is K in the K-means method. A third parameter is K, which is required for the K nearest neighbors (K n) method, but the author suggests using $K = \sqrt{n}$ as a default value.

The time complexity of the AVF, SCF, WAVF, and WDOD is $\approx O(nq)$, which linearly increases with the number of observations and the number of categorical variables.

The time complexity of the CMBOD method can be computed as follows: In each step, it scans the dataset once and it has 3 steps (3 scans). Consequently, the time complexity of CMBOD method is $\approx O(nq)$, which linearly increases with the number of observations and the number of categorical variables.

The time complexity of Bouguessa's method is composed of two parts: Computing outlying scores and identification of anomalies. It requires O(nq) in computing outlying scores. On the other hand, the complexity of the K-means method by using Lloyd's Algorithm [120] is O(nkqi), where i is number of iterations. While, the complexity of EM algorithm relies on the number of iterations and time for processing the expectation and maximization steps [29]. So, it approximately requires $O(im(t_E+t_M))$, where m is number of parameters to be estimated, t_E is the time required to do the expectation step, and t_M is the time required for the maximization step of the EM algorithm. Consequently, the total complexity of Bouguessa's method is $\approx O(nq + \text{Max}_m(nkqi + im(t_E + t_M)))$. Accordingly, it requires complex time processing. Another serious disadvantage of Bouguessa's method is that its output (identified outliers) is random, that is, the identified outliers can change from one run to the next due to the random nature of the K-means method.

The AVF, SCF, WAVF, WDOD, and CMBOD are fast and scalable methods (they can efficiently deal with large scale categorical datasets), but they only consider the marginal frequency and ignore any dependency structure among the categorical variables. Moreover, they require specifying the number of anomalies M in advance, which is impractical in real applications. In addition to the difficulty of defining the suitable value of these parameters, the results are very sensitive to these parameters.

3.2 Itemset Frequency

The methods in this category consider the frequency of itemsets (combinations of categories of size \leq *maxlen*, a predefined parameter). The anomalies are defined as observations that have small itemset frequencies. These methods start with building a set of frequent itemsets (itemsets with frequency \geq a certain parameter *minfreq*), then declare observations with less frequent itemsets

as outliers. Examples of these methods are: (a) Frequent Pattern Outlier Factor (FPOF) [89], (b) Link-based Outlier and Anomaly Detection in Evolving Datasets (LOADED) [73] and [131], (c) Outlier Detection for Mixed Attribute Datasets (ODMAD) [105], and (d) Frequent Non-Derivable Itemsets – Outlier Detection (FNDI-OD) [107].

The FPOF, the categorical part of LOADED, and the categorical part of ODMAD define anomalies as observations that contain infrequent patterns in their itemsets. For defining the list of frequent itemsets (S), they use two parameters: minfreq (minimum frequency of frequent itemset) and maxlen (maximum number of categories in the itemset). In addition to itemsets parameters, they use another parameter for declaring anomalies (e.g., FPOF declares M observations with infrequent itemsets as outliers; LOADED and ODMAD require two extra parameters; the window size (win) and the maximum difference score (ΔS). For each observation x_i , they define the window of x_i (a set of observations of size win). They compute the average outlying score in its window, then calculate the difference between the average of window scores and x_i 's score. If the difference exceeds ΔS , then they declare x_i as an outlier.

A significant problem with the itemset-based methods is that the number of frequent itemsets tends to be huge. For each observation, the set of frequent items is scanned to identify frequent items belonging to this observation. Thus, they take a very long time in processing. To reduce the running time of itemset-based anomaly detection methods, FNDI-OD replaces the set of all frequent itemsets with a condensed and compacted version called Non-Derivable Itemsets (NDI) by removing frequent derivable itemsets that can be deduced with their frequency from non-derivable itemsets. NDI represents the minimal frequent itemsets cover that can be used to reproduce all frequent itemsets. For example, if we have 3 frequent itemsets a, ab and abc with the same frequency 10. Then, ab is derivable itemset because its frequency can be derived from a and abc, while a and abc are non-derivable. The size of NDI is usually very small compared to the size of its relevant set of frequent items. FNDI-OD scans NDI for each observation. As a result, it significantly reduces the processing time compared to other itemset-based anomaly detection methods. At the same time, the FNDI-OD has similar accuracy to the FPOF [107].

The time-complexity of itemset frequency-based methods can be divided into two parts; Finding frequent itemsets and calculating outlying scores. An apriori method is used for finding the frequent itemsets [6]. Its time complexity is

$$O(n(C_1^Q + C_2^{Q-u} + C_3^{Q-u} + \dots + C_{maxlen}^{Q-u})) \approx O(n(2^{(Q-u)})),$$
(13)

where Q is the total number of categories and u is the number of infrequent categories [62]. In the worst case, all categories are frequent, then u=0. Accordingly, the complexity of the Apriori method in the worst case is $\approx O(n(2^Q))$. The time complexity of calculating outlying scores is O(n|S|). The total time complexity for these methods is then $O(n(2^Q + |S|))$. The FNDI-OD reduces the processing time of itemset-based anomaly detection methods, but in the worst case it has the same time complexity.

The itemset frequency-based methods take into account dependency among categorical variables, since they handle categorical variables jointly by considering itemset frequencies instead of only marginal frequencies. However, they have the following disadvantages: They are computationally expensive as they have exponential complexity with respect to the number of categorical variables and linear complexity with respect to the size of the set of frequent itemsets, which grows combinatorially. Moreover, they require two parameters; *maxlen* and *minfreq*. Their results are very sensitive to the values of these parameters.

3.3 Diversified Frequency

Under this category, we discuss one method for identifying anomalies in categorical data, that is, Couple Biased Random Walk (CBRW) [133]. The CBRW method takes into account intra-feature coupling (distribution) among categories at the same categorical variable as well as inter-feature coupling (interactions) among categories at different categorical variables. First, to begin with, let m_j be the modal (the most frequent) category for the categorical variable X_j and $p(m_j)$ be the relative frequency of the modal category of the categorical variable X_j . Then, for each category x_{ij} in X_j , we compute an intra-feature coupling deviation score as

$$\delta(x_{ij}) = [dev(x_{ij}) + base(m_i)]/2, \tag{14}$$

where $dev(x_{ij}) = [p(m_j) - p(x_{ij})]/p(m_j)$ and $base(m_j) = 1 - p(m_j)$. Thus, the intra-feature coupling considers both the marginal frequencies of categories as well as the whole frequency distribution within each variable.

Second, a directed graph, \mathcal{G} , is built to capture the inter-feature value coupling (interactions) among categorical values. In this graph, each category is represented as a node, while inter-feature coupling is propagated from a node u to another node v through an edge connecting u and v,

$$A(u,v) = p(u|v) = \frac{p(u,v)}{p(v)}.$$
(15)

The inter-feature coupling checks wether a certain category coupled with other outlying categories from other categorical variables. A(u, v) measures the strength of coupling between u and v. Finally, based on $\delta(v)$ is as defined in (14), and A(u, v), CBRW builds a biased random walk matrix W_b as

$$W_b(u,v) = \frac{\delta(v)A(u,v)}{\sum\limits_{v \in V} \delta(v)A(u,v)},$$
(16)

where V is the set of nodes in \mathcal{G} (all categories in the whole dataset). $W_b(u, v)$ represents the transition from the node u to the node v having a probability proportional to $\delta(v)A(u, v)$. Thus, each random walk is biased by the value of $\delta(v)$.

After building W_b , CBRW generates the probability distribution of biased random walk column vector, π_0 , to calculate an outlying score for each category. Then CBRW initializes π_0 values by a uniform distribution. Thereafter, it computes π_{t+1} at time t+1 on the basis of π_t as

$$\pi_{t+1} = (1 - \alpha) \frac{1}{|V|} 1 + \alpha W_b^t \pi_t, \tag{17}$$

where α is a damping factor to guarantee convergence. It is shown in [133] that if \mathcal{G} is irreducible and aperiodic, then π converges to a unique stationary probability, π^* , regardless of the initialization values π_0 . Subsequently, the outlying score for a certain category v is computed as

$$CBRW(v) = \pi^*(v). \tag{18}$$

The computed outlying CBRW scores can be used in features selection for high dimensional data and/or identifying outliers. CBRW scores can measure the importance of a certain categorical variable as

$$Rel(X_j) = \sum_{i=1}^{c_j} CBRW(x_{ij}), \tag{19}$$

where (x_{ij}) is the i^{th} category in the j^{th} variable and c_j is number of categories in the j^{th} variable. Thus, it can be utilized in ranking and selection of the most relevant categorical variables then it could be embedded into any outliers detection method.

Moreover, CBRW scores can identify outliers by computing an outlying score for each observations (x_i) as

$$CBRW-OD(x_i) = \sum_{j=1}^{q} W_j * CBRW(x_{ij}),$$
(20)

where

$$W_{j} = \frac{Rel(X_{j})}{\sum\limits_{i=1}^{q} Rel(X_{j})}.$$
(21)

Accordingly, CBRW labels observations having the highest M CBRW-OD scores as outliers.

The time complexity of the CBRW method is composed of three parts. First, building the biased random walk matrix W_b requires $O(nQ^2)$, where Q is total number of categories in all variables. Second, the complexity of computing the probability distribution of biased random walk vector is $O(|E|I_{max})$, where |E| is number of edges in \mathcal{G} and I_{max} is maximum number of iterations for calculating π^* . Finally, the complexity of finding the outliers is O(nq). Therefore, the overall complexity of CBRW is $O(nQ^2 + |E|I_{max} + nq)$.

One advantage of the CBRW is that it takes into consideration correlation among the categorical variables, since it captures frequency distribution of categorical variables as well as inter-feature coupling. In addition, it ranks categorical variables according to its relevance to outlierness. However, it is computationally expensive as it has a quadratic complexity with respect to the number of categories. Moreover, it requires two parameters: the number of outliers M and the damping factor α .

4 BAYESIAN/CONDITIONAL FREQUENCY-BASED METHODS

This approach has a different definition of categorical anomalies. It searches for observations with high marginal frequencies and small joint frequencies. In other words, it defines anomalies to be observations with infrequent combination of categories while its categories are frequent. The conditional anomalies are quite different from anomalies identified by other anomaly detection methods. They look for observations that contain frequent categories which are rarely observed together. Examples of methods that follow this approach are (a) Conditional Algorithm (CA) [50], (b) the Anomaly Pattern Detection (APD), (c) the method proposed in [144] (denoted here by RHH, the first letters of the authors' names), and (d) the Attribute Association (AA) algorithm [127].

The Conditional Algorithm (CA) defines outliers as observations with infrequent combinations of frequent categories, [50]. The CA calculates a measure of rareness, named $r(x_{ij}, x_{ik})$ ratio, between a pair of categories x_{ij} and x_{ik} in the categorical variables X_j and X_k , respectively. The $r(x_{ij}, x_{ik})$ ratio is defined as

$$r(x_{ij}, x_{ik}) = \frac{p(x_{ij}, x_{ik})}{p(x_{ij})p(x_{ik})},$$
(22)

where $p(x_{ij})$ and $p(x_{ij}, x_{ik})$ are the marginal and joint relative frequencies, respectively. A small r-ratio means that the multiplication of marginal probabilities is much higher than their joint probability, which indicates a higher probability of anomalous co-occurrence. The r-ratio is applied to itemsets, a subset of categorical attributes, up to maxlen length. For each observation, r-ratio is computed for every pair of itemsets. Then CA assigns the smallest r-ratio to be its outlying score, which means the outlying score reflects a subset of variables rather than considering all variables.

The CA method labels observations that have outlying scores less than a certain threshold θ as anomalies. The number of itemsets is huge because it grows combinatorially with the number of categorical variables, q. Therefore, the conditional algorithm tries to decrease the number of

candidate subsets by pruning itemsets from weakly correlated variables. It employs the mutual information ([158])

$$\mu(X_j, X_l) = -\sum_{i=1}^{i=c_j} \sum_{k=1}^{k=c_l} p(X_j = x_{ij} \land X_l = x_{kl}) \times \log \left(\frac{p(X_j = x_{ij} \land X_l = x_{kl})}{p(X_j = x_{ij}) \times p(X_l = x_{kl})} \right), \tag{23}$$

which measures the correlation between a pair of categorical variables X_j and X_l . It only considers itemsets within strongly correlated variables whose mutual information, μ , is greater than predefined threshold, β_{μ} . Consequently, it dramatically decreases the number of r-ratio computations.

The Anomaly Pattern Detection (APD) method is proposed to detect anomalous patterns, groups of related observations having outliers percentage higher than expected [51]. The APD method uses the CA algorithm [50] as first step to identify individual outliers. Then it uses rule-based technique to study the behavior of outliers in each pattern.

Similar to the conditional probability method, the RHH method [144] searches for observations with frequent attribute values but infrequent joint co-occurrence. It can find anomalies in categorical datasets as well as mixed dataset. In the mixed datasets, it transforms quantitative variables into categorical variables by discretizing quantitative variables into fixed length intervals. In the training phase, it builds Bayesian network (see, e.g., [37] and [126]) to capture the dependency among the attributes. It takes advantage of the parent-child relationship of Bayesian network and stores the counts of categories and combinations of categories using the Alternating Decision (AD) Tree [125]. It computes the same outlying score as the conditional algorithm.

The Attribute Association (AA) algorithm [127] defines conditional anomalies as observations that contain frequent categories but these itemsets are rarely observed together. The method starts with deriving a set of association rules with high confidence from the data. Then it computes an outlying score called outlier degree as

$$od(x_i) = \frac{|C_{x_i}^+ - C_{x_i}^0|}{|C_{x_i}^0|},$$
(24)

where $C_{x_i}^0$ and $C_{x_i}^+$ are initial and final cover sets for x_i . For details about how to construct these cover sets, see [127]. The method labels an observation x_i as an outlier if $od(x_i) > \theta$, where θ is a predefined threshold.

The time complexity of the CA, the RHH, and the outlier detection part of the APD can be divided into two parts: Finding frequent itemsets and calculating the r-ratio by building the Bayesian network. The time complexity for finding frequent itemsets is $O(n(2^Q))$, where Q is the total number of categories. However, the complexity of building Bayesian network is $\approx O(q^w)$, where w is the tree width [47]. Therefore, the total complexity of the conditional algorithm and the RHH method is $\approx O(n(2^Q) + (q^w))$.

On the other hand, the time complexity for the AA algorithm consists of three parts: Finding frequent itemsets, identifying a set of association rules with high confidence, and building the final cover. The first and most expensive part is finding frequent itemsets. The time complexity for finding frequent itemsets is $\approx O(n((2^Q)))$. The complexity of identifying confident association rules is $\binom{|f|}{2}$, where |f| is the length of frequent itemsets, which is combinatorial in the number of frequent itemsets f. The time complexity for finding the final cover is $(|R|)^2$, where |R| is the length of the set of high confidence association rules. Therefore, the total complexity of the AA method is $\approx O(n((2^Q))) + (C_2^{|f|}) + (|R|)^2)$.

Conditional anomalies detection methods have time complexity problems, since they require a combinatorial time for building itemsets and a long time for searching in the high conditional probability space. Moreover, they require many parameters (see, e.g., minfreq, maxlen, and minconf).

5 DENSITY-BASED METHODS

Density-based anomalies or local anomalies approach aims at identifying observations that have outlying behavior in local areas, in which observations usually share similar characteristics [35]. Local anomalies are different from global anomalies which are inconsistent with the pattern suggested by the majority of all other observations, not only observations within their local areas [35], [182], [44], and [98]. Local anomaly detection methods for categorical data include (a) the Hyperedge-based Outlier Test (HOT) [176], (b) the k-Local Anomalies Factor k-LOF [182], and (c) the WATCH method [112].

The HOT method is based on defining two sets of variables: Common variables, C, and Outlying variables, A. HOT ignores other variables. It develops a graph called hypergraph that captures the similarity among categorical observations. The hypergraph is divided into groups (named hyperedges, H) based on the frequent itemsets in C. Each hyperedge $he \in H$ contains a set of observations that has the corresponding itemset he. Then, for each observation x_i in the hyperedge he and an attribute A_j in a set of outlying attributes A, HOT calculates a deviation measure,

$$Dev^{he}(x_i, A_j) = \frac{S_{A_j}^{he}(c_r) - \mu_{S_{A_j}^{he}}}{\sigma_{S_{A_j}^{he}}},$$
(25)

where c_r is the value of x_i in the jth categorical variable and $S_{A_j}^{he}(x_i)$ is the number of observations in hyperedge he which have the value c_r in the categorical attribute A_j ,

$$\mu_{S_{A_j}^{he}} = \frac{1}{|A_j^{he}|} \sum_{c_r \in A_j^{he}} S_{A_j}^{he}(c_r)$$
 (26)

is the mean value of $S_{A_i}^{he}(c_r)$, and

$$\sigma_{S_{A_j}^{he}} = \sqrt{\frac{1}{|A_j^{he}|} \sum_{c_r \in A_j^{he}} (S_{A_j}^{he} - \mu_{S_{A_j}^{he}})^2}$$
 (27)

is the standard deviation of $S_{A_j}^{he}(c_r)$. HOT declares an observation x_i as an outlier for the hyperedge he and outlying attributes A_i if $Dev^{he}(x_i, A_i) < \theta$, where θ is a predefined deviation parameter.

The k-LOF is a local anomaly detection method for both categorical and quantitative data [182]. It extends Local Anomalies Factor (LOF) [35] to categorical domain. The k-LOF identifies an observation as a local outlier if its relationships with its neighbors are weaker than the relationships between its neighbors and its neighbors' neighbors. It measures the similarity among observations in the categorical data by a weighted undirected graph named similarity graph, G = (V, E, w), where V is a set of vertices (observations) in the graph G, E is a set of edges (unordered pairs of vertices), and w is the weight of similarity between two vertices.

The k-LOF uses the concept of k-walk, which is the connection of length k edges between two observations. $N^k(x_i)$ is the set of neighbors connected to x_i within k-walks including x_i . The k-LOF defines the similarity of k-walk between two observations x_i and x_j , as

$$s^{k}(x_{i}, x_{j}) = \begin{cases} w(x_{i}, x_{j}), & \text{if } (x_{i}, x_{j}) \in E \text{ and } k = 1, \\ \sum_{(x_{m}, x_{j}) \in E} w(x_{m}, x_{j}) \times s^{k-1}(x_{i}, x_{m}), & \text{if } k > 1. \end{cases}$$
(28)

In addition, the k-LOF computes another type of similarity named the accumulated similarity of k-walk between two observations x_i and x_j , as

$$S^{k}(x_{i}, x_{j}) = \sum_{i=1}^{k} s^{i}(x_{i}, x_{j}).$$
(29)

Then, the above two similarity measures are combined in an outlying score for x_i , which is given by

$$k\text{-LOF}(x_i) = \frac{1}{S^k(x_i, x_i) \times n} \sum_{j=0}^n S^k(x_i, x_j).$$
 (30)

The k-LOF labels an observation, x_i , as an outlier if k-LOF(x_i) > θ , where θ is a predefined parameter. The k-LOF takes into consideration the direct relationships between an observation and its direct neighbors as well as the indirect relationships among the neighbors and the neighbors' neighbors where k > 1. Thus, it requires the parameters θ and the maximum length of indirect relationships k.

Most recently, the WATCH method for identifying outliers in high dimensional categorical datasets using feature grouping is proposed in [112]. The method aims at detecting local outliers (outliers in subspaces of the full-dimensional space). WATCH consists of two phases: Feature grouping and outliers detection. In the feature grouping phase, it partitions the set of q categorical variables into q disjoint feature groups by collecting correlated variables into the same group. Let

$$FR(X_j, X_l) = \frac{I(X_j, X_l)}{\mathcal{E}(X_i, X_l)}$$
(31)

be the feature relation between two features X_i and X_l , where

$$I(X_j, X_l) = -\sum_{i=1}^{l=c_j} \sum_{k=1}^{k=c_l} p(X_j = x_{ij} \land X_l = x_{kl}) \times \log \left(\frac{p(X_j = x_{ij} \land X_l = x_{kl})}{p(X_j = x_{ij}) \times p(X_l = x_{kl})} \right), \tag{32}$$

is the Mutual Information [158], which measures the correlation between a pair of categorical variables X_j and X_l , and

$$\mathcal{E}(X_j, X_l) = -\sum_{i=1}^{i=c_j} \sum_{k=1}^{k=c_l} p(X_j = x_{ij} \land X_l = x_{kl}) \times \log(p(X_j = x_{ij} \land X_l = x_{kl}))$$
(33)

is the Shannon Entropy [160] between a pair of categorical variables X_i and X_l .

WATCH performs the following steps, which are similar to the K-means clustering algorithm [67], to cluster the categorical variables into $g \in 2, 3, \cdots, q/2$ feature groups: Initialization of the pivot dimensions, allocating each feature to one of the feature groups, and updating the selection of the pivot dimensions. During the initialization of the pivot dimensions, a random feature is selected to be the pivot dimension, η_1 , for the first feature group, \mathcal{G}_1 , then the next pivot dimension $\eta_i, i \in 2, 3, \cdots, g$, is the feature that has the minimum sum of feature relations with all the selected pivot dimensions, $\sum_{k=1}^{i-1} \mathrm{FR}(\eta_i, \eta_k)$.

In the allocation step, for each categorical variable, X_j , WATCH searches for the optimal feature group \mathcal{G}_i that has the maximum feature relation with its pivot dimension, $FR(\eta_i, X_j)$. Then, it updates the choice of the pivot dimensions in each feature group according to the allocation step. Finally, it iterates the allocation of each categorical attribute to the optimal (best) feature group and updating the pivot dimensions until there is no change in the pivot dimensions and hence, there is no change in feature groups.

For any specific partition of $q \in 2, 3, \dots, q/2$ groups, the aggregate relation,

$$AR(D) = \sum_{i=1}^{i=g} \sum_{r \in \mathcal{G}_i} FR(\eta_i, X_r),$$
(34)

is computed for the whole data set. WATCH chooses the value of g that maximizes the aggregate relation, that is,

$$g = \operatorname{argmax}_{q \in 2, 3, \dots, q/2} \operatorname{AR}(D). \tag{35}$$

The second phase of WATCH method is detecting outliers in each feature group. It calculates a weighting factor, $w(X_j)$ for each categorical variable according to the feature correlation between this variable and other variables within its feature group,

$$w(X_j) = -\frac{1}{r} \sum_{l=1}^r FR(X_l, X_j),$$
(36)

where r is the number of categorical variables in the feature group and $FR(X_l, X_j)$ is the feature relation between X_l and X_j . Then, it computes an outlying score for each observation according to each feature group G_k as

WATCH_{$$G_k$$} $(x_i) = \frac{1}{r} \sum_{j=1}^r \begin{cases} 0, & \text{if } f(x_{ij}) = 1, \\ w(X_j) \times \log\left(\frac{(f(x_{ij})-1)^{(f(x_{ij})-1)}}{f(x_{ij})^{f(x_{ij})}}\right), & \text{otherwise.} \end{cases}$ (37)

The outlying score, WATCH $\mathcal{G}_k(x_i)$, takes negative values. Furthermore, the higher outlying score an observation has, the higher probability of being an outlier. Thus, WATCH labels observations having the highest M scores as outliers with respect to each feature group. It takes the union of outliers sets in all feature groups to declare all outlying observations. Therefore, the number of declared outliers by WATCH is at most $M \times g$.

Accordingly, the WATCH method takes into consideration correlation among categorical attributes by partitioning dimensions into g feature groups. The method aims at detecting hidden outliers (those observations that significantly differ from other observations in subsets of correlated dimensions).

The time complexity of HOT is $\approx O(nq^2 + q|H| max(|he|))$, where |H| is the total number of hyperedges (frequent itemsets) and max(|he|) is the maximum number of observations found in a single hyperedge. The HOT has a linear complexity with respect to the number of observations and the size of the hyper graph, while, it has a quadratic complexity with respect to the number of categorical variables.

The time complexity of k-LOF can be divided into two parts: Building similarity graph G and computing similarity and outlying scores. The time complexity for building the similarity graph is $O(n^2q)$. However, the time complexity of computing similarity and outlying scores is $\approx O(n^2\sum_{i=1}^k i)$. Therefore, the total time complexity of k-LOF is $\approx O(n^2(q+\sum_{i=1}^k i))$.

The time complexity of the WATCH algorithm consists of three parts. First, the building feature groups phase requires O(qgt), where t is number of iterations. Second, the complexity of detecting outliers in all feature groups requires O(nsg). Third, the complexity of finding the optimal number of feature groups is O((q/2-1)qgt). Therefore, the overall complexity of WATCH is $O(q^2gt/2+ngs)$.

Accordingly, the density-based anomaly detection methods, HOT, k-LOF, and WATCH are computationally expensive due to the quadratic cost of measuring similarity with respect to the number of observations. Additionally, these methods require specifying decision parameters k and θ in advance. Specifying suitable values of these parameters is difficult and the results are also very sensitive to these parameters.

6 CLUSTERING-BASED METHODS

Clustering-based anomaly detection methods rely on clustering categorical datasets then identify observations that are located in sparse regions as outliers. There are two methods that belong to this category: (a) Ranking-based Outliers Analysis and Detection (ROAD) [167] and [169], and (b) the Rough-ROAD method [168] and [170].

The Ranking-based Outliers Analysis and Detection (ROAD) method defines two types of outliers; frequency-based and clustering-based outliers. Frequency-based outliers are those observations which have infrequent categories (small average marginal frequencies). However, clustering-based outliers are those observations which have infrequent combinations of frequent categories.

ROAD develops two different ranking schemes; one for each outlier-type. First, it computes a density score for each observation (the average marginal frequencies) as $den(x_i) = \frac{1}{q} \sum_{j=1}^q f(x_{ij})$, which is the same as AVF(x_i) in equation(1). ROAD sorts observations based on their density scores. Then it gives the higher probability for being frequency-based outliers to those observations having small density scores.

Second, ROAD partitions the given categorical dataset into k clusters using the k-mode algorithm [92]. Then it defines the set of big clusters, BC, as the clusters that contain at least $\alpha\%$ of observations, where α is the big cluster threshold. Additionally, ROAD computes the distance between an observation x_i and a certain cluster C_l , $d(x_i, C_l)$, as $d(x_i, C_l) = \sum_{j=1}^q \phi(x_{ij}, z_{lj})$, where z_{lj} is the j^{th} category in the representative (mode), z_l , in the l^{th} cluster, C_l , and $\phi(x_{ij}, z_{lj})$ is computed as

$$\phi(x_{ij}, z_{lj}) = \begin{cases} 1, & \text{if } x_{ij} \neq z_{lj}, \\ 1 - \frac{|C_{lj}|}{|C_l|}, & \text{otherwise,} \end{cases}$$
 (38)

where $|C_{lj}|$ is the number of observations that have x_{ij} and $|C_l|$ is the number of observations in the l^{th} cluster.

ROAD defines a cluster-based ranking scheme, $clus-rank(x_i)$ based on the distance to the nearest big cluster. The larger distance to nearest big cluster, the higher probability for being cluster-based outliers. ROAD defines two sets of most likely outliers by labeling top M observations in each rank. For the sake of integration, ROAD defines an observation as an outlier if it belongs to frequency-based outliers or cluster-based outliers. It computes the union of the sets of outliers. Accordingly, the ROAD takes into account marginal frequencies as well as joint frequencies. However, it requires two parameters; M and the big cluster threshold, α .

The Rough-ROAD method extends the ROAD method. Here, the k-modes algorithm [92] is modified to the rough k-modes algorithm based on using rough sets theory [136] to capture the uncertainty and deal with ambiguity regarding the membership of outliers to certain clusters. Rough-ROAD consists of two phases; rough clustering phase and ranking outliers phase. The only difference between ROAD and Rough-ROAD is in the clustering phase.

Similar to the original k-modes algorithm, the rough k-modes aims at partitioning categorical data into k clusters. Each cluster, C_j , has a representative called C_j 's mode, Z_j , the observation with the highest density in C_j . The rough k-modes consists of three basic steps: the modes random initialization step, the assigning of each observation to clusters, and the updating of the modes based on the assigning step. Rough k-modes and k-modes are iterative algorithms repeating the second and the third steps until convergence (no change in clusters' distributions).

For each cluster, C_j , rough k-modes defines two sets of observations a lower approximation, $\underline{C_j}$, and an upper approximation, $\overline{C_j}$. The lower approximation $\underline{C_j}$ contains all objects that are surely located in cluster C_j . In other words, it contains all objects, where C_j is the only close cluster to these objects. Whereas, the upper approximation $\overline{C_j}$ contains all objects that are close to two or

more clusters. After the mode initialization step, the Rough-ROAD method assigns an object to either a single lower approximation set of a rough cluster or multiple upper approximation sets of rough clusters. For each object, x_i , it identifies the nearest cluster, C_j , which has the minimum $d(x_i, Z_j)$, where $1 \le j \le k$. Then, it looks for the set of clusters, T, that are close with respect to its nearest cluster, that is,

$$T = \left\{ C_l : \frac{d(x_i, Z_l)}{d(x_i, Z_i)} \le \varepsilon \right\},\tag{39}$$

where $\varepsilon \ge 1$ is the roughness parameter. If the set of close clusters T, is empty, then x_i has only one close rough cluster C_j . Thus, x_i is assigned only to the lower approximation of the j^{th} cluster, $\underline{C_j}$. However, if T is not empty, then x_i has two or more close rough clusters C_j and C_l . Thus, x_i is assigned to the upper approximation of those clusters $\overline{C_j}$ and $\overline{C_l}$.

After Rough-ROAD assigns all objects to rough clusters, it recalculates the new clusters' modes, $Z_1, Z_2, ..., Z_k$. In each rough cluster, C_j , it searches for the object that has the maximum density to be the mode, Z_j . Then it iterates assigning objects to rough clusters and updating clusters' modes until convergence (there is no change in the clusters).

After the rough clusters are built, outliers detection in Rough-ROAD is the same as the ROAD method. It is based on building a ranking scheme for each type of outliers.

Accordingly, the Rough-ROAD is based on the rough set theory to capture the uncertainty and deal with ambiguity regarding the dataset taking into account the density of objects within its cluster as well as the marginal frequencies. However, it requires the following 5 parameters: The number of outliers M, the roughness parameter ε , the lower approximation weight parameter w_l , the number of clusters k, and the big cluster threshold α .

The time complexity of ROAD consists of three parts. First, building the frequency-based ranking scheme requires $O(nqc_{max})$, where c_{max} is maximum number of categories per categorical variable. Second, the complexity of computing the cluster-based ranking scheme needs $O(nqk^2 + nqkt)$, where k is the number of clusters and t is the number of iterations required for convergence in the k-mode algorithm. Third, the complexity of ranking observation is $O(n \log n)$. Therefore, the overall complexity of ROAD is $O(nqc_{max} + nqk^2 + nqkt + n \log n)$. The time complexity of Road-ROAD is as complex as that of ROAD. Therefore, the clustering-based methods are computationally expensive.

7 DISTANCE-BASED METHODS

The distance-based anomaly detection approach for categorical data extends the concept of distance-based anomalies for quantitative data ([104], [103], [142] and [16]). There are many definitions of anomalies in distance-based approach. These include:

- Anomalies are the M observations whose average distances to the k nearest neighbors are the greatest [17] and [16].
- Anomalies are the *M* observations whose distances to the *k*-th nearest neighbor are the greatest [142].
- Anomalies are the observations that have fewer than p observations within a certain distance d [104] and [103].
- Anomalies are observations that have the highest z-scores of the average distances to the k-nearest neighbors. That is, first we compute the average distance of each observation to its k-nearest neighbors. Then, we standardized these average distances and obtain the z-scores of the average distances. The observations with z-scores greater than a threshold θ (e.g., 3) are declared as outliers [63]. This method does not require the parameter M in advance, but

it assumes that the z-scores follow a standard normal distribution to help in choosing value of θ .

Examples of distance-based methods for categorical data are (a) Orca (name of software) [24], (b) a method called iOrca [28], (c) the Common Neighbor Based distance (CNB) [113], and (d) the Recursive Binning and Re-Projection (RBRP) method [74].

Orca identifies anomalies in mixed datasets. The categorical part of Orca uses the Hamming distance to compute the distance between two categorical observations. The Hamming distance between two categorical observations is the number of mismatches between them [93]. The output of Hamming distance is an integer value between 0 and q, the number of categorical variables in the data. For each observation x_i , Orca computes an outlying score which is the average Hamming distance between x_i and its k nearest neighbors. Orca declares k observations with the highest outlying scores as outliers.

On the other hand, a method called indexed Orca (iOrca), [28], is proposed to speed up Orca. iOrca depends on a simple index technique by choosing a random observation, x_r , most likely to be inlier, and then rank all observations based on the distance from x_r . Then, observations are assessed based on their distances to x_r . iOrca can be up to an order of magnitude faster than Orca but in the worst case (when the randomly chosen reference observation is an outlier) iOrca becomes worse than Orca.

The CNB is similar to Orca but it relies on similarity-based distances [40] instead of the Hamming distances. The similarity-based distance between a pair of objects x_i and x_m , $d_{\text{CNB}}(x_i, x_m, \theta)$, is computed in steps. First, finding the common neighbors set

$$CNS(x_i, x_m, \theta) = NS(x_i) \cap NS(x_m), \tag{40}$$

where $NS(x_i)$ is the set of nearest neighbors for (x_i) , which contains an object, x_r , that has $sim(x_i, x_r) \ge \theta$, where $sim(x_i, x_r)$ is the number of matched categories and θ is a predetermined parameter. Second, the distance between x_i and x_m , $d_{CNB}(x_i, x_m, \theta)$, is computed as

$$d_{\text{CNB}}(x_i, x_m, \theta) = 1 - \left(\frac{\log_2 | \text{CNS}(x_i, x_m, \theta) |}{\log_2 n}\right). \tag{41}$$

CNB declares the M observations whose average distances to the k nearest neighbors are the highest as anomalies.

The time complexity of Orca is $O(n^2q)$ in the worst case, whereas the time complexity of CNB is $O(n^2(k + S(\theta) + q) + n(k + M))$, where $S(\theta)$ is the average number of neighbors.

The distance-based anomaly detection methods are not efficient for very large datasets due to its quadratic time complexity with respect to the number of observations [24], [74], and [28]. Orca performs a pruning strategy to improve the performance of the distance-based anomaly detection methods. For each observation, x_i , if its score during the calculation is less than a certain threshold (the smallest score of anomalies in that iteration), then further processing on that observation is not needed because it can no longer be an outlier. In the worst case (when observations are arranged in ascending order according to Orca scores), Orca still takes a quadratic time. However, the authors mentioned that Orca runs close to linear time in practice [74].

The Recursive Binning and Re-Projection (RBRP) method is proposed in [74] to reduce the time complexity of distance-based anomaly detection methods to $O(nq \log n + qn^2)$, where n is number of observations and q is number of variables. RBRP groups observations in clusters and then, for each observation x_i , it searches for its nearest neighbors in its cluster first then the next closest cluster and so on. In practice, RBRP outperforms Orca by an order of magnitude, [74].

Distance-based anomaly detection methods have the following problems. Their results are very sensitive to the values of the parameters θ and k (number of nearest neighbors). Moreover, they

require identifying the number of anomalies M in advance, which is impractical in real applications. They ignore the dependency among categorical variables as well as the number of categories c_j . Finally, they have expensive time complexity $O(n^2q)$ in the worst case. In addition, the outputs of RBRP and iOrca are random because the initial observation is chosen at random in the case of iOrca and the RBRP uses the K-mode algorithm, the output of which is random.

8 INFORMATION THEORETIC METHODS

The idea behind information theoretic approach relies on the direct relationship between the existence of anomalies and the amount of noise in the dataset. Consequently, the problem of anomaly detection can be transformed into an optimization problem, where the objective is finding the set of anomalies that maximizes the information gain, or minimizing the uncertainty, of the inlier observations. Most of these methods use an entropy as an information gain or as an uncertainty measure. These methods can be further classified into 2 groups based on the type of entropy they use: (a) the Shannon entropy and (b) the Holo entropy.

8.1 Methods Based on the Shannon Entropy

The Shannon entropy of a dataset X is defined as [160]:

$$\mathcal{E}(X) = -\sum_{x_i \in X} p(x_i) \log p(x_i), \tag{42}$$

where x_i is an observation in X and $p(x_i)$ is the relative frequency of x_i in X. This can be interpreted as the amount of information contained in the dataset.

The methods based on the Shannon entropy search for a set of observations that minimize the Shannon entropy of the remaining observations. Examples of these methods are (a) the Local Search heuristic Algorithm (LSA) [87] and (b) the Greedy Algorithm (GA) [88].

The LSA finds the M observations, thought to be outliers, that minimize the Shannon entropy of $X_{(I)}$ (the dataset X without the M observations indexed by I). The LSA algorithm builds all C_M^n possible combinations of M, computes the entropy $\mathcal{E}(X_{(I)})$ for each combination, and chooses I that gives the smallest value of $\mathcal{E}(X_{(I)})$. The time complexity of LSA is $O(nqC_M^n)$. LSA algorithm scans the dataset C_M^n times. Therefore, it is infeasible for large n and M.

Since LSA algorithm is not applicable for large data due to its huge time complexity, the GA [88] is proposed to decrease the complexity of LSA. It assumes that the observations are independent from each other. The GA searches for the observation that minimizes the entropy of the remaining (n-1) observations, then it declares that observation as an anomaly and removes it from X. It repeats for the remaining observations until M anomalies have been declared. The time complexity of the GA is O(nqM). Therefore, it significantly decreases the number of dataset scans from C_M^n to M. Although, the GA substantially decreases time complexity of LSA algorithm, its time complexity is still high.

The Shannon entropy-based methods, the LSA and the GA, have the following problems. The time complexity is too high especially when the number of observations and the number of anomalies are large. Moreover, they require identifying the number of anomalies in advance, M, which is impractical in real applications. They are based on the joint frequency which is always very low especially in large datasets (large number of categorical variables). Therefore, the problem of ties happens due to many subsets giving the same minimum value of the entropy. In addition, they are subject to the masking (failing to declare some outliers) and swamping (declaring some inliers as outliers) problems.

8.2 Methods Based on the Holo Entropy

The total correlation,

$$C_t(X) = \sum_{i=1}^{q} \mathcal{E}(X_i) - \mathcal{E}(X), \tag{43}$$

measures the shared information within a dataset. It is the sum of mutual information of all variables in a dataset X. It can be interpreted as a measure of purity or goodness of the dataset. While the Holo entropy is defined as

$$HE(X) = \mathcal{E}(X) + C_t(X) = \mathcal{E}(X) + \sum_{j=1}^{q} \mathcal{E}(X_j) - \mathcal{E}(X) = \sum_{j=1}^{q} \mathcal{E}(X_j).$$
 (44)

Consequently, the total correlation as well as the Shannon entropy are aggregated together into one measure to compute information gain in the dataset.

The Holo entropy-based anomaly detection methods search for the M observations, thought to be outliers, that minimize the Holo entropy of $X_{(I)}$ (the dataset X without M observations indexed by I) [179], [180], and [140]. To avoid multiple scans for each observation, x_i , (as in the GA), the differential Holo entropy is calculated, $h^x(x_i)$, as the change in the Holo entropy between the dataset with and without x_i . The higher differential Holo entropy an observation has, the higher probability of being an outlier. Then the M observations with the highest differential Holo entropy are labeled as outliers. Examples of the Holo entropy-based algorithms are (a) the Information Theoretic-Based (ITB) [180] and (b) the Excess Entropy-Based (EEB) [152]. Each of these can use either a single pass (SP) or a step by step (SS). Accordingly, we have four algorithms: ITB-SP, ITB-SS, EEB-SP, and EEB-SS.

In the single pass the dataset is scanned only once then observations with the highest outlier factor (differential entropy) are declared as outliers. The step by step means scanning data sets M times. In each time, the observation with the highest outlier factor (differential Holo entropy) is declared as an outlier and we repeat until M observations are declared as outliers.

The difference between ITB and EEB is in the type of entropy they use. The ITB relies on the Holo Entropy, whereas, the EEB is based on computing excess entropy \mathcal{E}_E , also known as the dual total correlation or binding information. The excess entropy quantifies the amount of entropy present in a dataset after subtracting the sum of the entropies of each variable conditioned upon all other variables. It is calculated as

$$\mathcal{E}_{E}(X) = \mathcal{E}(X) - \sum_{j=1}^{q} \mathcal{E}(X_{j} | (X \setminus X_{i})), \tag{45}$$

where $\mathcal{E}(X \setminus X_i)$ is the dataset X with the variable X_i removed.

The time complexity of ITB-SS and EEB-SS is O(nqM) but the time complexity of ITB-SP and EEB-SP is O(nq). Thus, they significantly decrease the number of dataset scans from C_m^n to either one dataset scan in the case of ITB-SP and EEB-SP or to M dataset scans in the case of ITB-SS and EEB-SS. The Holo entropy-based methods, ITB and EEB require identifying the number of anomalies in advance, M. The ITB-SP and EEB-SP scan dataset only once, which means they are fast and scalable but they may have the masking and swamping problems.

9 COMPRESSION-BASED METHODS

Compression algorithms are usually used in the fields of communication and storage rather than in data mining. Recently, however, few compression-based anomaly detection methods are proposed based on the fact that anomalies do not comply with the model suggested by the other points in

the data. Accordingly, observations that could not compress well are considered as outliers. These methods look for the best compression model that suits the non-outlying data points. Objects that deviate (have bad compression measures) are highlighted as anomalies [31]. Examples of compression-based anomaly detection methods for categorical data are (a) KRIMP [164] and (b) Comprex [15].

KRIMP makes use of the Minimum Description Length (MDL) compression concept to decide whether a categorical observation was drawn from the training distribution or it will be considered as an outlier. MDL finds the optimal compression model, \ddot{m} , from the set of compression models M that minimizes the length of both the encoded dataset and the compression model. An observation, x_i , is highlighted as an anomaly if $L(x_i|\ddot{m}) > \theta$, where $L(x_i|\ddot{m})$ is the encoded length of x_i by the optimal compression model, \ddot{m} , and θ is a decision parameter.

KRIMP is based on the idea of code tables. A code table consists of two columns. The first column contains the itemsets and the second column contains their codes. Itemsets are sorted in descending order according to their lengths then their frequencies. Higher order (that is, longer and more frequent) itemsets take shorter code. Each observation is represented by a set of non-overlapping itemsets that completely cover all vales of that observation. An anomaly can be seen as an observation which contains infrequent itemsets and hence, it's encoded code is longer than other observations.

Comprex is similar to KRIMP but it uses multiple code tables instead of one. It makes use of correlation among categorical attributes in building a dictionary (code table) for each strongly correlated set of attributes. It measures the variables correlation by the Information Gain (IG). IG is measured as the average number of bits we can save when compressing X_i and X_j together. If the saving equals zero, it means that no saving from considering X_i and X_j jointly, hence they are independent. The goal is to find the optimal set of dictionaries that could achieve the minimal lossless compression for the given dataset. Comprex declares observations that have long encoded length as outliers. The idea behind using multiple code tables is to improve compression measures. Moreover, Comprex builds code tables directly from the data by using strongly associated variables instead of using itemset. Therefore, it dose not require building itemsets and hence it avoids its sensitive parameters and expensive cost. The authors claim that Comprex is a parameter-free method, but it actually requires a decision parameter, M, the number of anomalies, which have the highest compression cost. Comprex makes use of Cantelli's Inequality [77] to estimate an approximate value of M.

The time complexity of KRIMP is divided into two parts: Finding frequent itemsets and encoding the full dataset. The complexity of finding frequent itemsets by apriori algorithm [6] is $\approx O(n(2^Q))$ and the time complexity of encoding the full dataset is $f^2(nQ+1)$, where n is the number of observations, f is the number of frequent itemsets and Q is the total number of categories. Thus, the total time complexity for KRIMP is $\approx O(n(2^Q+f^2(nQ+1)))$. The time complexity of Comprex is $O(nq^2)$, which has a quadratic complexity with respect to the number of categorical variables.

Although compression-based anomaly detection methods for categorical data outperform other relevant methods in terms of error and detection rates, they are computationally very expensive especially for datasets that contain large number of variables. Similar to other anomaly detection methods for categorical data, compression-based algorithms require decision parameters to decide whether an observation is or is not an outlier.

10 SEMISUPERVISED METHODS

Anomaly identification methods can be classified on the basis of the availability of training samples into three approaches: Supervised, Unsupervised, and Semisupervised [91]. Supervised anomaly detection methods may be considered as a binary classification problem, where each observation is

categorized into one of two known classes: normal (inliers) and abnormal (anomalies). Supervised anomaly detection methods consist of two phases: The training phase and the testing phase. In the training phase, well-known samples from both classes, normal and abnormal, are required in building a classification model to distinguish between the behavior of anomalous observations and normal ones. The established classification model is then used for predicting the class label for the observations in the testing phase [137].

The unsupervised anomaly detection methods are used when prior information about the class label are not available. These methods try to divide the data into two unbalanced groups, normal and abnormal, based on the belief that normal observations are located in dense regions while abnormal observations are located in sparse regions [35].

The semisupervised anomaly detection methods use samples of normal observations to build a model of normal data in the training phase. Hence, they do not require abnormal samples in the training phase. Then, they classify instances which significantly deviate from the normal model as anomalies [91].

Supervised techniques are usually faster and more accurate than unsupervised techniques but they require prior information about normal and abnormal observations which enhances their performances. However, in some applications, the prior information especially for abnormal data are not available or expensive to obtain (e.g., nuclear fault detection). Moreover, different types of anomalies, which did not appear in the training phase, may appear in the testing phase. As a result, semisupervised anomaly detection approach is useful [90].

A commonly used semisupervised anomaly detection methods for quantitative data is the Oneclass Support Vector Machine (OSVM) [157]. Examples of semisupervised anomaly detection for categorical data are Feature Regression and Classification (FraC) [128] and [129] and a Semisupervised Anomaly Detection framework for Categorical data (SAnDCat) [95].

FRaC is a semisupervised feature modeling approach for detecting anomalies in quantitative and categorical data. During the training phase, FRaC makes use of supervised learning algorithms to build ensemble predictive models. In the testing phase, it labels observations that disagree with the predictive models as anomalies. Initially, the set of features (attributes) should be identified for which predictive models will be learnt. A predictive model uses a predefined set of features to predict the values of a certain feature. Next, a supervised learning algorithm is chosen to train feature models such as regression models for quantitative attributes and decision trees for categorical attributes. FRaC may use multiple supervised learning models to predict one feature then it combines them into a feature ensemble model to improve learning accuracy.

In the testing phase, an anomaly score for each observation is computed based on all feature predictors (e.g., the average number of correct predictions per observation). A higher score indicates stronger agreement with the predictive models and a lower score indicates higher probability of being an anomaly. The preceding anomaly score does not take into account uncertainty of classifier feature models. Therefore, FRaC computes the surprisal anomaly score for each observation x_i as

$$FRaC(x_i|C) = \sum_{i=1}^{D} surprisal(P(x_{ij}|C_j, \rho_i(x_i))),$$
 (46)

where C is the set of predictive models, D is number of features for which predictive models are learnt, and $(C_j, \rho_j(x_i))$ is the predicted value of the j-th attribute in x_j based on the mapping function $\rho_i(x_i)$ of the predictor C_j . The surprisal anomaly score is interpreted as the amount of evidence that an instance is anomalous [128].

SAnDCat is also a semisupervised anomaly detection framework for categorical data based on Distance Learning for Categorical Attributes (DILCA) [94], which calculates the normalized

Euclidean distance between the conditional probability of categories, (y_i, y_j) , given the values of other attributes called the context attributes X_k . DILCA computes a context-based distance between two categorical values, y_i and y_j , based on the similarity between their probability distribution using the contingency table [94].

SAnDCat consists of two steps. At the beginning, it learns from the training data to build a normality model, \mathcal{M} , consisting of one matrix for each categorical attribute, X_j . The elements of this matrix represent the DILCA distance between each pair of categories in the corresponding attribute. Then, SAnDCat computes anomaly score for each instance in the testing data, x_i , as the average distance between x_i and selected k representative instances from the training set. SAnDCat framework includes four heuristics to select k representatives; random, central k instances in the data, closest k instances for each observation, and farthest k instances for each observation. Finally, the SAnDCat either labels k observations having the highest anomaly score as outliers or labels observations with scores exceeding a certain threshold, k, as outliers [94].

The complexity and accuracy of FRaC is sensitive to the selected predictors which are used to calculate feature models. Some predictors are linear with respect to sample size (e.g., Naive Bayes). Moreover, FRaC executes multiple cross-validation loops for each feature and for each classifier of the ensemble classifiers for each feature. Therefore, the complexity of FRaC in the worst case scenario is $O(n^2q)$ [95]. The time complexity of SAnDCat can be divided into three parts: Building the classification model, choosing k representatives, and making decision. The complexity of building classification model in the training phase using DILCA distance metric is $O(nq^2 \log q)$ [94]. The complexity of choosing the k representatives from the training data in the worst case is $O(n^2)$. The complexity of making decision in the worst case is $O(n \log n)$. Consequently the time complexity of SAnDCat is $O(nq^2 \log q + n^2 + n \log n)$.

Semisupervised anomaly detection methods for categorical data outperform other methods in terms of detection rate and computational time. Unlike other anomaly detection approaches for categorical data, the semisupervised approaches require normal instances to learn the normal behavior in the training phase.

11 DISCUSSIONS

Existing anomaly detection methods for categorical data face many problems. Three important problems are discussed below:

11.1 Computational Complexity

Time computational complexity is an important issue in data mining research since most of real applications have huge datasets in terms of the number of observations and the number of categorical variables. Accordingly, the time complexity is an important criterion for choosing anomaly detection methods for categorical data. The time complexities for various methods are given in Table 1 in the Appendix.

The methods FPOF, LOADED, ODMAD, FNDI-OD, AA, CA and KRIMP are based on finding frequent itemsets. They require an exponential time complexity with respect to the number of categorical variables and several passes over the input dataset. Consequently, these techniques are not recommended for datasets with large number of categorical variables.

Distance-based and density-based methods (Orca, CNB, RBRP, iOrca, *k*-LOF, Indicator Variable, MCA, FRaC, SAnDCat, and MCA using Canberra distance) always take a very long time for distance computations due to nested loops. Measuring pair-wise distances requires a quadratic time complexity with respect to number of observations. Therefore, these techniques are not recommended for datasets with large number of observations.

The time complexities of the methods based on the Shannon entropy, LSA and GA, are $O(nqC_M^n)$ and O(nqM), respectively, and hence LSA is infeasible for large datasets. The GA, ITB-SS, and EEB-SS scan the dataset M times. Consequently, they take very long time especially when the number of observations and number of anomalies are large.

The time complexities of Comprex and CBRW-OD are approximately $O(nq^2)$ and $O(nQ^2 + |E|I_{max} + nq)$, respectively. Therefore, these techniques are not recommended for datasets with large number of categorical variables. The time complexity of ROAD and Rough-ROAD is based on k^2 , where k is the number of clusters. The time complexities of ROAD, Rough-ROAD and WATCH are affected by the number of iterations until convergence (finding good or optimal solution).

The time complexity of AVF, SCF, ITB-SP, EEB-SP, WAVF, WDOD, and CMBOD is O(nq). It linearly increases with respect to the number of observations and the number of categorical variables. Thus they are efficient in handling huge datasets, but they ignore dependency among categorical variables.

11.2 Human Intervention

Human intervention (input parameters) is another problem for anomaly detection methods for categorical data. Existing methods require one or more input parameters. In real applications, defining the suitable values of these parameters is a hard and critical task. In addition, the results are sensitive to these input parameters. The most commonly required parameters are given in Table 2 in the Appendix.

- (1) Decision parameters: Existing anomaly detection methods for categorical data compute a score for each observation. To identify whether an observation is an outlier, they require one of the following parameters:
 - *M*: The number of assumed anomalies in the dataset. Most of existing methods (AVF, SCF, CMBOD, FPOF, FNDI-OD, CBRW-OD, WATCH, ROAD, Rough-ROAD, CNB, Orca, CNB, LSA, GA, ITB-SP, EEB-SP, ITB-SS, EEB-SS and CompreX) require *M* in advance.
 - θ: Outlying score threshold. Observations with scores greater than θ are declared as outliers. The methods that use θ are the Indicator Variables method, MCA, WDOD, CA, RHH, AA, HOT, k-LOF, KRIMP, FRaC, and SAnDCaT.
 - win and $\triangle S$: If the difference between an observation score and the average of observations scores in its win is greater than $\triangle S$, then that observation is declared as an outlier. The parameters win and $\triangle S$ are used in LOADED and ODMAD.
- (2) Itemset-based parameters: (*minfreq* and *maxlen*) are used in defining frequent itemsets. They are required by FPOF, LOADED, ODMAD, FNDI-OD, HOT, AA, CA and KRIMP.
- (3) *k*: The number of nearest neighbors in distance-based anomaly detection methods for categorical data. It is required in Bouguessa's method, Orca, CNB, and *k*-LOF.

11.3 Real and Synthetic Categorical Datasets

Testing the performance of anomaly detection methods requires both real benchmark data as well as synthetic data. Real benchmark quantitative data are abundant. By contrast, there are only few benchmark categorical datasets available such as Breast Cancer and Lymphography datasets, [68].

In addition, generating quantitative data with and without anomalies is straight forward because the definition of anomalies in quantitative data is simple, since anomalies are defined to be a minority of observations that are inconsistent with the pattern suggested by the majority of observations in a dataset. Accordingly, quantitative data with anomalies are generated as follows: The inlier observations are generated from a statistical distribution with certain mean and standard deviation, whereas anomalies are generated from another statistical distribution with different mean and/or

standard deviation. The inliers and outliers are then merged together and used for testing the performance of methods for the detection of anomalies in quantitative data. Computation times and detection performance measures (detection rate, precision and recall) can be computed and used to compare various methods using the same datasets.

By contrast, generating synthetic categorical data and planting anomalies inside them is more difficult than it is for the case for quantitative dataset [36]. The main reason is the lack of agreement in the literature about the definition of outliers. Some anomaly detection papers use synthetic categorical datasets but only to test time performance not detection performance. One way to overcome this difficulty is to first generate quantitative data with outliers, then discretize (categorize) the variables to obtain synthetic categorical data which can be used to compute time performance and detection performance measures [171].

12 CONCLUSIONS AND FUTURE WORK

Anomaly detection is very important and useful in numerous real-life applications such as identifying computer network intrusions and fraud detection. We surveyed available methods for the identification of anomalies for categorical data in the statistics as well as the machine learning and computer science literatures. There is no general agreement on a single definition of an anomaly in categorical data. We reviewed 36 methods for anomaly detection in categorical data in both the statistics and computer science literatures and classified them into 12 different categories based on the conceptual definition of anomalies they use. We identified the strength and weakness of each method. The computational complexity as well as the required parameters for each method are provided, since they are critical problems in real applications. In addition, we have discussed the common challenges in the detection of anomalies in categorical data.

There are several directions for further research in anomaly detection in categorical data. Extending anomaly detection for categorical data to new research areas such as recommendation systems [56], categorical data streams [39], [124], moving objects [71], [116], and information network [5], [13] and [69] may be studied in future research. Defining an automatic critical values instead of specifying the number of anomalies in advance is another direction for future research.

Table 1. Various Anomaly Detection Algorithms and their Complexities.

Approach	Method	Complexity
Indicator Variables	IV (•)	$O(nq + n^2 \sum_{j=1}^{q} (c_j - 1))$
indicator variables	MCA	$O(nq + 2n^2Q + nQ^2 + Q^3)$
	AVF	O(nq)
	SCF	O(nq)
Manginal Enggyanay	WAVF	$\approx O(nq)$
Marginal Frequency	WDOD	$\approx O(nq)$
	CMBOD	$\approx O(nq)$
	Bouguessa's method (*) (•)	$\approx O(nq + \text{Max}_m(nkqi + im(t_E + t_M)))$
	FPOF (•)	$\approx O(n(2^Q))$
Itemset Frequency	LOADED (*) (•)	$\approx O(n(2^Q))$
itemset Frequency	ODMAD (*) (•)	$\approx O(n(2^Q))$
	FNDI-OD (●)	$\approx O(n(2^Q))$
Diversified Frequency	CBRW-OD (●)	$O(nQ^2 + E I_{max} + nq)$
	CA (•)	$\approx O(n(2^Q) + (q^w))$
Payasian/Conditional Engage	APD (•)	$\approx O(n(2^Q) + (q^w))$
Bayesian/Conditional Frequecy	RHH (*) (●)	$\approx O(n(2^Q) + (q^w))$
	AA	$\approx O(n((2^Q)) + (C_2^{ f }) + (R)^2)$
	HOT	$\approx O(nq^2 + H q \max(he))$
Density-Based Methods	k-LOF (*)	$O(n^2(q+\sum_{i=1}^k i))$
	WATCH (*)	$O(ngs + q^2gt/2)$
Chartenia a Deced Methods	ROAD (●)	$O(nqc_{max} + nqk^2 + nqkt + n\log n)$
Clustering-Based Methods	Rough-ROAD (●)	$O(nqc_{max} + nqk^2 + nqkt + n \log n)$
	Orca (*)	$O(qn^2)$ linear in practice
Distance-Based Methods	CNB	$O(n^2(k+S(\theta)+q)+n(k+M))$
Distance-based Methods	RBRP	$O(nq\log n + qn^2)$
	iOrca	$O(qn^2)$
Shannon Entropy	LSA (●)	$O(nqC_M^n)$
Shaimon Entropy	GA	O(nqM)
	ITB-SP	O(nq)
Hala Entuany	EEB-SP	O(nq)
Holo Entropy	EEB-SS	O(Mnq)
	EEB-SS	O(Mnq)
Compression	KRIMP	$\simeq O(n(2^Q + f^2(nQ+1))))$
Compression	CompreX	$\approx O(nq^2)$
SemiSupervised	FraC (*)	$\approx O(n^2q)$
-	SAnDCat	$\simeq O(nq^2\log q + n^2 + n\log n)$

^(*) This complexity is for the categorical part of methods for mixed data.

^(•) This complexity is derived and reported in this article.

Table 2. Various Anomaly Detection Algorithms and their Required Parameters.

			No	Decision	Intrinsic
Approach	Algorithm	Reference	Param.	Parameters	Parameters
Indicator Variables	IV	[162]	1	θ	
mucator variables	MCA	[162]	1	θ	
	AVF	[106]	1	W	
	SCF	[173]	1	W	
Monginal Descentation	WAVF	[150]	1	W	
marginar riequency	WDOD	[186]	1	θ	
	CMBOD	[109]	1	W	
	Bouguessa's (*)	[34]	3	$M_{ m max}, knn$	knn
	FPOF	[88]	3	W	minfreq, maxlen
Itemset Frequency	LOADED (*)	[131]	4	win, △S	minfreq, maxlen
	ODMAD (*)	[105]	4	win, △S	minfreq, maxlen
	FNDI-OD	[107]	3	W	minfreq, maxlen
Diversified Frequency	CBRW-OD	[133]	2	W	α
	CA	[20]	2	θ	minfreq, maxlen, $lpha$, eta_lpha
Ravecian/Conditional Erecusess	APD	[51]	2	θ	minfreq, maxlen, $lpha$, eta_lpha
payesiam committee inchact	RHH (*)	[144]	4	θ	minfreq, maxlen, $lpha$
	AA	[127]	4	θ	minfreq, maxlen, minconf
	HOT	[176]	3	θ	minfreq, maxlen
Density-Based	k-TOF (*)	[182]	2	θ	k-walks
	WATCH	[112]	1	W	
Clustering-Based	ROAD	[169]	3	W	α, K
Ciusici ing-Daseu	Rough-ROAD	[170]	2	W	$\varepsilon, w_l, K, \alpha$
	Orca (*)	[24]	2	W	knn
Distance-Based	CNB	[113]	3	M	knn, sim
Distance Dasca	RBRP	[74]	3	W	knn, K
	iOrca	[28]	2	W	knn
Channon Entrony	TSA	[87]	1	W	
Snamon Enuopy	GA	[88]	1	W	
	ITB-SP	[180]	1	W	
Holo Entrony	EEB-SP	[152]	1	W	
tion run ob)	ITB-SS	[180]	1	W	
	EEB-SS	[152]	1	W	
Compression	KRIMP	[164]	3	θ	minfreq, maxlen
Compression	CompreX	[15]	1	W	
SemiSupervised	FraC (*)	[129]	1	θ	
	SAnDCat	[95]	2	θ	k

(*) The parameters required for the categorical part of methods for mixed data.

REFERENCES

- [1] Abduvaliyev, A., Pathan, A.-S. K., Zhou, J., Roman, R., and Wong, W.-C. (2013). On the vital areas of intrusion detection systems in wireless sensor networks. *IEEE Communications Surveys & Tutorials*, 15(3):1223–1237.
- [2] Abukhalaf, H., Wang, J., and Zhang, S. (2015). Outlier detection techniques for localization in wireless sensor networks: A survey. *International Journal of Future Generation Communication and Networking*, 8(6):99–114.
- [3] Aggarwal, C. C. (2017). Outlier Analysis. 2nd ed. Springer, Cham.
- [4] Aggarwal, C. C. and Yu, P. S. (2001). Outlier detection for high dimensional data. In *Proceedings of the ACM International Conference on Management of Data, SIGMOD*, pages 37–46, Santa Barbara, California, USA.
- [5] Aggarwal, C. C., Zhao, Y., and Yu, P. S. (2011). Outlier detection in graph streams. In *Proceedings of the ACM IEEE International Conference on Data Engineering, ICDE*, pages 399–409, Hannover, Germany.
- [6] Agrawal, R. and Srikant, R. (1994). Fast algorithms for mining association rules in large databases. In *Proceedings of International Conference on Very Large Data Bases, VLDB*, pages 487–499, Santiago, Chile.
- [7] Agresti, A. (2010). Analysis of Ordinal Categorical Data. John Wiley & Sons, New York, NY, USA, 2nd edition.
- [8] Agresti, A. (2013). Categorical Data Analysis. John Wiley & Sons, New York, NY, USA, 3rd edition.
- [9] Agyemang, M., Barker, K., and Alhajj, R. (2006). A comprehensive survey of numeric and symbolic outlier mining techniques. *Intelligent Data Analysis*, 10(6):521–538.
- [10] Ahmed, M., Mahmood, A. N., and Hu, J. (2016a). A survey of network anomaly detection techniques. *Network and Computer Applications*, 60:19–31.
- [11] Ahmed, M., Mahmood, A. N., and Islam, M. R. (2016b). A survey of anomaly detection techniques in financial domain. *Future Generation Computer Systems*, 55:278–288.
- [12] Ajitha, P. and Chandra, E. (2015). A survey on outliers detection in distributed data mining for big data. *Journal of Basic and Applied Scientific Research*, 5(2):31–38.
- [13] Akoglu, L., Mcglohon, M., and Faloutsos, C. (2010). Oddball: Spotting anomalies in weighted graphs. In *Proceedings of the Pacific Asia Knowledge Discovery and Data Mining, PAKDD*, pages 420–431, Hyderabad, India.
- [14] Akoglu, L., Tong, H., and Koutra, D. (2015). Graph based anomaly detection and description: A survey. *Data Mining and Knowledge Discovery*, 29(3):626–688.
- [15] Akoglu, L., Tong, H., Vreeken, J., and Faloutsos, C. (2012). Fast and reliable anomaly detection in categorical data. In *Proceedings of the ACM International Conference on Information and Knowledge Management, CIKM*, pages 415–424, HI, USA
- [16] Angiulli, F., Basta, S., and Pizzuti, C. (2006). Distance-based detection and prediction of outliers. *IEEE Transaction on Knowledge and Data Engineering*, 18(2):145–160.
- [17] Angiulli, F. and Fassetti, F. (2002). Fast outlier detection in high dimensional spaces. In *Proceedings of the European Conference on the Principles of Data Mining and Knowledge Discovery*, pages 19–26, Helsinki, Finland.
- [18] Ankur, Y. N. and Singh, A. S. (2014). Oulier analysis using frequent pattern mining: A review. *International Journal of Computer Science and Information Technologies*, 5(1):47–50.
- [19] Archana, N. and Pawar, S. (2014). Survey on outlier pattern detection techniques for time-series data. *International Journal of Science and Research (IJSR)*, 1(1):1852–1856.
- [20] Bailetti, T., Gad, M., and Shah, A. (2016). Intrusion learning: An overview of an emergent discipline. *Technology Innovation Management Review*, 6(2):15–20.
- [21] Bakar, U. A. B. U. A., Ghayvat, H., Hasanm, S. F., and Mukhopadhyay, S. C. (2016). Activity and anomaly detection in smart home: A survey. In Mukhopadhyay, S. C., editor, *Next Generation Sensors and Systems*, chapter 9, pages 191–220. Springer, New York, NY, USA.
- [22] Bakar, Z. A., Mohemad, R., Ahmad, A., and Deris, M. M. (2006). A comparative study for outlier detection techniques in data mining. In *Proceedings of IEEE International Conference on Cybernetics and Intelligent Systems*, pages 1–6, Bangkok, Thailand.
- [23] Barnett, V. and Lewis, T. (1994). Outliers in Statistical Data. John Wiley & Sons, New York, NY, USA, 3rd edition.
- [24] Bay, S. and Schwabacher, M. (2003). Mining distance-based outliers in near linear time with randomization and a simple pruning rule. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 29–38, Washington, DC, USA.
- [25] Beh, E. J. (2008). Simple correspondence analysis of nominal-ordinal contingency tables. *Journal of Applied Mathematics and Decision Sciences*, 228:1–17.
- [26] Beldar, A. P. and Wadne, V. S. (2015). The detail survey of anomaly/outlier detection methods in data mining. *International Journal of Multidisciplinary and Current Research*, 3:462–472.
- [27] Bezerra, C. G., Costa, B. S. J., Guedes, L. A., and Angelov, P. P. (2015). A comparative study of autonomous learning outlier detection methods applied to fault detection. In *IEEE International Conference on Fuzzy Systems*, *FUZZ-IEEE*, pages 1–7, Istanbul, Turkey.

- [28] Bhaduri, K., Matthews, B. L., and Giannella, C. R. (2011). Algorithms for speeding up distance-based outlier detection. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 895–867, San Diego, CA, USA.
- [29] Bhagyashree, U. and Nilav, M. (2014). Overview of k-means and expectation maximization algorithm for document clustering. In *International Conference on Quality Up-gradation in Engineering, Science and Technology (ICQUEST)*, pages 5–8, Maharashtra, India.
- [30] Billor, N., Hadi, A. S., and Velleman, P. (2000). Blocked adaptive computationally-efficient outlier nominators. *Computational Statistics and Data Analysis*, 34:279–298.
- [31] Böhm, C., Haegler, K., Müller, N. S., and Plant, C. (2009). Coco: coding cost for parameter-free outlier detection. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 149–158, Paris, France.
- [32] Boriah, S., Chandola, V., and Kumar, V. (2008). Similarity measures for categorical data: A comparative evaluation. In *Proceedings of the International SIAM Data Mining Conference, SDM*, pages 243–254, Atlanta, GA, USA.
- [33] Bouguessa, M. (2014). A mixture model-based combination approach for outlier detection. *International Journal on Artificial Intelligence Tools*, 23(4):1–21.
- [34] Bouguessa, M. (2015). A practical outlier detection approach for mixed-attribute data. Expert Systems with Applications, 42:8637åÅ\$–8649.
- [35] Breunig, M. M., Kriegel, H., Ng, R. T., and Sander, J. (2000). LOF: Identifying density-based local outliers. In *Proceedings of the ACM International Conference on Management of Data, SIGMOD*, pages 93–104, Dallas, Texas, USA.
- [36] Campos, G. O., Zimek, A., Sander, J., Campello, R. J., Micenková, B., Schubert, E., Assent, I., and Houle, M. E. (2016). On the evaluation of unsupervised outlier detection: Measures, datasets, and an empirical study. *Data Mining and Knowledge Discovery*, 30(4):891–927.
- [37] Castillo, E., Gutirrez, J. M., and Hadi, A. S. (1997). Expert Systems and Probabilistic Network Models. Springer-Verlag, New York, NY, USA.
- [38] Chandola, V., Banerjee, A., and Kumar, V. (2009a). Anomaly detection: A survey. ACM Computing Surveys (CSUR), 41(3):1–58.
- [39] Chandola, V., Banerjee, A., and Kumar, V. (2012). Anomaly detection for discrete sequences: A survey. *Transactions on Knowledge and Data Engineering*, 24(5):823 839.
- [40] Chandola, V., Boriah, S., and Kumar, V. (2008). Understanding categorical similarity measures for outlier detection. Technical report, University of Minnesota, Department of Computer Science and Engineering,1-46.
- [41] Chandola, V., Boriah, S., and Kumar, V. (2009b). A framework for exploring categorical data. In *Proceedings of the International SIAM Data Mining Conference, SDM*, pages 187–198, Sparks, NV.
- [42] Chatterjee, S. and Hadi, A. S. (1986). Influential observations, high leverage points, and outliers in regression. *Statistical Sciences*, 1:379–416.
- [43] Chatterjee, S. and Hadi, A. S. (1988). Sensitivity Analysis in Linear Regression. John Wiley & Sons, New York, NY, USA.
- [44] Chawla, S. and Sun, P. (2006). Slom: A new measure for local spatial outliers. *Knowledge and Information Systems*, 9:412–429.
- [45] Cheng, H., Tan, P.-N., Potter, C., and Klooster, S. A. (2009). Detection and characterization of anomalies in multivariate time series. In *Proceedings of the SIAM International Conference on Data Mining*, SDM, pages 413–424, Lisbon, Portugal.
- [46] Cho, H. and Eo, S.-H. (2016). Outlier detection for mass spectrometric data. In Jung, K., editor, *Statistical Analysis in Proteomics*, chapter 5, pages 91–102. Springer, New York, NY, USA.
- [47] Cooper, G. F. (1990). The computational complexity of probabilistic inference using bayesian belief networks. Artificial Intelligence, $42:393-\exists E_1405$.
- [48] Cousineau, D. and Chartier, S. (2015). Outliers detection and treatment: A review. *International Journal of Psychological Research*, 3(1):58–67.
- [49] Daniel, J. V., Joshna, S., and Manjula, P. (2013). A survey of various intrusion detection techniques in wireless sensor networks. *International Journal of Computer Science and Mobile Computing*, 2(9):235–246.
- [50] Das, K. and Schneider, J. (2007). Detecting anomalous records in categorical datasets. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 220–229, San Jose, CA, USA.
- [51] Das, K., Schneider, J., and Neill, D. B. (2008). Anomaly pattern detection in categorical datasets. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 169–176, Las Vegas, NV, USA.
- [52] Dave, D. and Varma, T. (2014). A review of various statistical methods for outlier detection. *International Journal of Computer Science & Engineering Technology (IJCSET)*, 5(2):137–140.
- [53] Debar, H., Dacier, M., and Wespi, A. (1999). Towards a taxonomy of intrusion-detection systems. *Computer Networks*, 31(9):805–822.
- [54] D'Enza, A. I. and Greenacre, M. (2012). Multiple correspondence analysis for the quantification and visualization of large categorical data sets. In Di Ciaccio, A., Coli, M., and Ibañez, J. M. A., editors, Advanced Statistical Methods for the

- Analysis of Large Data-Sets, pages 453-463. Springer.
- [55] Deshmukh, M. M. K. and Kapse, A. (2016). A survey on outlier detection technique in streaming data using data clustering approach. *International Journal of Engineering and Computer Science*, 5(1):15453–15456.
- [56] Desrosiers, C. and Karypis, G. (2011). A comprehensive survey of neighborhood-based recommendation methods. In *Recommender Systems Handbook*, pages 107–144. Springer-Verlag New York, NY, USA.
- [57] Devi, R. L. and Amalraj, R. (2015). Hubness in unsupervised outlier detection techniques for high dimensional data-a survey. *International Journal of Computer Applications Technology and Research*, 4(11):797–801.
- [58] Dhimmar, J. H. and Chauhan, R. (2014). A survey on profile-injection attacks in recommender systems using outlier analysis. *International Journal of Advance Research in Computer Science and Management Studies*, 2(12):356–359.
- [59] Ding, X., Li, Y., Belatreche, A., and Maguire, L. P. (2014). An experimental evaluation of novelty detection methods. *Neurocomputing*, 135:313–327.
- [60] Divya, K. T. and Kumaran, N. S. (2016). Survey on outlier detection techniques using categorical data. *International Research Journal of Engineering and Technology*, 3:899–904.
- [61] Dokas, P., Ertoz, L., Kumar, V., Lazarevic, A., Srivastava, J., and Tan, P.-N. (2002). Data mining for network intrusion detection. In the Proceedings of the NSF Workshop on Next Generation Data Mining, pages 21–30, Baltimore, MD, USA.
- [62] Du, J., Zheng, Q., Li, H., and Yuan, W. (2005). The research of mining association rules between personality and behavior of learner under web-based learning environment. In *Proceedings of the the International Conference on Advances in Web-Based Learning ICWL 2005*, pages 15–26, Hong Kong, China.
- [63] Ebdon, D. (1991). Statistics in Geography: A Practical Approach-Revised with 17 Programs. Wiley-Blackwell, Hoboken, NI, USA.
- [64] Emran, S. M. and Ye, N. (2001). Robustness of canberra metric in computer intrusion detection. In *Proceedings of the IEEE Workshop on Information Assurance and Security*, pages 80–84, New York, NY, USA.
- [65] Fanaee-T, H. and Gama, J. (2016). Tensor-based anomaly detection: An interdisciplinary survey. *Knowledge-Based Systems*, 98:130–147.
- [66] Faria, E. R., Gonalves, I. J. C. R., de Carvalho, A. C. P. L. F., and Gama, J. (2015). Novelty detection in data streams. *Artificial Intelligence Review*, 45(2):235–269.
- [67] Forgy, E. W. (1965). Cluster analysis of multivariate data: efficiency versus interpretability of classifications. *Biometrics*, 21:768ï£i-780.
- [68] Frank, A. and Asuncion, A. (2018). UCI machine learning repository, http://archive.ics.uci.edu/ml/datasets.html, March, 2018
- [69] Gao, J., Liang, F., Fan, W., Wang, C., Sun, Y., and Han, J. (2010). On community outliers and their efficient detection in information networks. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 813–822, Washington DC, USA.
- [70] Garcia-Teodoro, P., Diaz-Verdejo, J., Maciá-Fernández, G., and Vázquez, E. (2009). Anomaly-based network intrusion detection: Techniques, systems and challenges. *computers & security*, 28(1):18–28.
- [71] Ge, Y., Xiong, H., Zhou, Z.-H., Ozdemir, H., Yu, J., and Lee, K. (2010). Top-eye: Top-k evolving trajectory outlier detection. In *Proceedings of the ACM Conference on Information and Knowledge Management, CIKM*, pages 1–4, Toronto, Canada.
- [72] Ghosh, D. and Vogt, A. (2012). Outliers: An evaluation of methodologies. In *Joint Statistical Meetings*, pages 3455–3460. American Statistical Association, San Diego, CA, USA.
- [73] Ghoting, A., Otey, M. E., and Parthasarathy, S. (2004). Loaded: Link-based outlier and anomaly detection in evolving data sets. In *Proceedings of the IEEE International Conference on Data Mining, ICDM*, pages 387–390, Ohio State, USA.
- [74] Ghoting, A., Parthasarathy, S., and Otey, M. E. (2008). Fast mining of distance-based outliers in high dimensional datasets. *Data Mining and Knowledge Discovery Journal, DMKD*, 16(3):349–364.
- [75] Gogoi, P., Bhattacharyya, D., Borah, B., and Kalita, J. K. (2011). A survey of outlier detection methods in network anomaly identification. *The Computer Journal*, 54(4):570–588.
- [76] Golub, G. H. and van Loan, C. F. (2012). Matrix computations, 3rd edition. John Hopkins U. Press.
- [77] Grimmett, G. and Stirzaker, D. (2001). Probability and Random Processes. 3rd ed. Oxford University Press, Oxford, UK.
- [78] Gunamani, V. and Abarna, M. (2013). A survey on intrusion detection using outlier detection techniques. *International Journal of Science, Engineering and Technology Research (IJSETR)*, 2(11):2063 –2068.
- [79] Gupta, M., Gao, J., Aggarwal, C. C., and Han, J. (2014a). Outlier detection for temporal data. Synthesis Lectures on Data Mining and Knowledge Discovery, 5(1):1–129.
- [80] Gupta, M., Gao, J., Aggarwal, C. C., and Han, J. (2014b). Outlier detection for temporal data: A survey. *IEEE Trans. Knowl. Data Eng.*, 26(9):2250–2267.
- [81] Hadi, A. S. (1992a). Identifying multiple outliers in multivariate data. *Journal of the Royal Statistical Society, Series (B)*, 54:761–771.

- [82] Hadi, A. S. (1992b). A new measure of overall potential influence in linear regression. *Computational Statistics and Data Analysis*, 14:1–27.
- [83] Hadi, A. S. (1994). A modification of a method for the detection of outliers in multivariate samples. *Journal of the Royal Statistical Society, Series (B)*, 56:393–396.
- [84] Hadi, A. S., Imon, A. H. M. R., and Werner, M. (2009). Detection of outliers. Wiley Interdisciplinary Reviews: Computational Statistics, 1:57–70.
- [85] Hadi, A. S. and Simonoff, J. S. (1993). Procedure for the identification of outliers in linear models. *Journal of the American Statistical Association*, 88:1264–1272.
- [86] Han, X., Yan, Y., Cheng, C., Chen, Y., and Zhu, Y. (2014). Monitoring of oxygen content in the flue gas at a coal-fired power plant using cloud modeling techniques. *IEEE Transactions on Instrumentation and Measurement*, 63(4):953–963.
- [87] He, Z., Xu, X., and Deng, S. (2005a). An optimization model for outlier detection in categorical data. In *Proceedings of the International Conference on Advances in Intelligent Computing*, pages 400–409, Hefei, China.
- [88] He, Z., Xu, X., and Deng, S. (2006). A fast greedy algorithm for outlier mining. In *Proceedings of the Pacific Asia Knowledge Discovery and Data Mining, PAKDD*, pages 567–576, Singapore.
- [89] He, Z., Xu, X., Huang, J. Z., and Deng, S. (2005b). Fp-outlier: Frequent pattern based outlier detection. *Computer Science and Information Systems, ComSIS*, 2:726–732.
- [90] Hido, S., Tsuboi, Y., Kashima, H., Sugiyama, M., and Kanamori, T. (2011). Statistical outlier detection using direct density ratio estimation. *Knowledge Information Systems*, 26(2):309âĂŞ-336.
- [91] Hodge, V. and Austin, J. (2004). A survey of outlier detection methodologies. Artificial Intelligence Review, 22:85-126.
- [92] Huang, Z. (1997). A fast clustering algorithm to cluster very large categorical data sets in data mining. In *Proceedings* of the International Data Mining and Knowledge Discovery, DMKM, Workshop at the ACM International Conference on Mangagement of Data, SIGKDD,, pages 1–8.
- [93] Huang, Z. and Ng, M. K. (1999). A fuzzy k-modes algorithm for clustering categoircal data. *IEEE Transaction and Fuzzy Systems*, 7:446–452.
- [94] Ienco, D., Pensa, R. G., and Meo, R. (2012). From context to distance: Learning dissimilarity for categorical data clustering. ACM Transactions on Knowledge Discovery from Data, 6(1):1–12.
- [95] Ienco, D., Pensa, R. G., and Meo, R. (2017). A semisupervised approach to the detection and characterization of outliers in categorical data. *IEEE Transactions on Neural Networks and Learning*, 28(5):1017–1029.
- [96] Ieva, F. and Paganoni, A. M. (2015). Detecting and visualizing outliers in provider profiling via funnel plots and mixed effect models. *Health Care Management Science*, 18(2):166–172.
- [97] Jiang, S., Song, X., Wang, H., Han, J.-J., and Li, Q.-H. (2006). A clustering-based method for unsupervised intrusion detections. *Pattern Recognition Letters*, 27:802–810.
- [98] Joshi, V. and Bhatnagar, R. (2014). Cbof: Cohesiveness-based outlier factor a novel definition of outlier-ness. In Proceedings of the International Workshop on Machine Learning and Data Mining in Pattern Recognition, MLDM, pages 175–189, Petersburg, Russia.
- [99] Joudaki, H., Rashidian, A., Minaei-Bidgoli, B., Mahmoodi, M., Geraili, B., Nasiri, M., and Arab, M. (2015). Using data mining to detect health care fraud and abuse: A review of literature. *Global Journal of Health Science*, 7(1):194–202.
- [100] Kalinichenko, L., Shanin, I., and Taraban, I. (2014). Methods for anomaly detection: A survey. In All-Russian Conference Digital Libraries: Advanced Methods and Technologies, Digital Collections, RCDL, pages 20–25, Dubna, Russia.
- [101] Kathiresan, V. and Vasanthi, N. A. (2015). A survey on outlier detection techniques useful for financial card fraud detection. *International Journal of Innovations in Engineering and Technology, IJIET*, 6(1):226–235.
- [102] Kaur, R. and Singh, S. (2015). A survey of data mining and social network analysis based anomaly detection techniques. *Egyptian Informatics Journal*, 39:1–18.
- [103] Knorr, E., Ng, R., and Tucakov, V. (2000). Distance-based outliers: Algorithms and applications. *VLDB Journal*, 8:237–253.
- [104] Knorr, E. M. and Ng, R. T. (1997). A unified approach for mining outliers. In *Proceedings of the International Conference of the Centre for Advanced Studies on Collaborative Research, CASCON*, pages 236–248, Toronto, Ontario, Canada.
- [105] Koufakou, A., Georgiopoulos, M., and Anagnostopoulos, G. (2008). Detecting outliers in high-dimensional datasets with mixed attributes. In *Proceedings of the International Conference on Data Mining, DMIN*, Los Vegas, NV, USA.
- [106] Koufakou, A., Ortiz, E., Georgiopoulos, M., Anagnostopoulos, G., and Reynolds, K. (2007). A scalable and efficient outlier detection strategy for categorical data. In *Proceedings of the IEEE International Conference on Tools with Artificial Intelligence, ICTAI*, pages 210–217, Patras-Peloponnese-Greece.
- [107] Koufakou, A., Secretan, J., and Georgiopoulos, M. (2011). Non-derivable itemsets for fast outlier detection in large high-dimensional categorical data. *Knowledge and Information Systems*, 29(3):697–725.
- [108] Lazarevic, A., Ertöz, L., Kumar, V., Ozgur, A., and Srivastava, J. (2003). A comparative study of anomaly detection schemes in network intrusion detection. In *Proceedings of the SIAM International Conference on Data Mining, SDM*, pages 25–36, San Francisco, CA, USA.

- [109] Lei, D., Zhang, L., and Zhang, L. (2013). Cloud model-based outlier detect algorithm for categorical data. *International Journal of Database Theory and Application*, 6(14):199–213.
- [110] Li, D. (2000). Uncertainty in knowledge representation. Chinese Engineering Science, 2(10):73-79.
- [111] Li, J. and Guo, J. (2015). A new feature extraction algorithm based on entropy cloud characteristics of communication signals. *Mathematical Problems in Engineering*, 2015:1–8.
- [112] Li, J., Zhang, J., Pang, N., and Qin, X. (2018). Weighted outlier detection of high-dimensional categorical data using feature grouping. *IEEE Transaction on Systems, Man, and Cybernetics: Systems*, pages 1–14.
- [113] Li, S., Lee, R., and Lang, S.-D. (2007). Mining distance-based outliers from categorical data. In *Proceedings of the IEEE International Conference on Data Mining Workshops, ICDM*, pages 225–230, Vancouver, Canada.
- [114] Liang, J. Y., Chin, K. S., and Dang, C. Y. (2002). A new method for measuring uncertainty and fuzziness in rough set theory. *Int J Gen Syst*, 31:331–342.
- [115] Lin, S. and Brown, D. E. (2006). An outlier-based data association method for linking criminal incidents. *Decision Support Systems*, 41:604–615.
- [116] Liu, W., Zheng, Y., Chawla, S., Yuan, J., and Xie, X. (2011). Discovering spatio-temporal causal interactions in traffic data streams. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD*, pages 1010–1018, San Diego, CA, USA.
- [117] Liu, X., Chen, F., and Lu, C.-T. (2014). On detecting spatial categorical outliers. GeoInformatica, 18(3):501-536.
- [118] Mahapatro, A. and Khilar, P. M. (2013). Fault diagnosis in wireless sensor networks: A survey. *IEEE Communications Surveys & Tutorials*, 15(4):2000–2026.
- [119] Malik, K., Sadawarti, H., and Kalra, G. S. (2014). Comparative analysis of outlier detection techniques. *International Journal of Computer Applications*, 97(8):12–21.
- [120] Manning, C. D., Raghavan, P., and Schutze, H. (2008). *Introduction to information retrieval*. Cambridge University Press, Cambridge, UK.
- [121] Marinho, J., Granjal, J., and Monteiro, E. (2015). A survey on security attacks and countermeasures with primary user detection in cognitive radio networks. *EURASIP Journal on Information Security*, 2015(1):1–14.
- [122] Markou, M. and Singh, S. (2003a). Novelty detection: A review-part 1: Statistical approaches. *Signal Processing*, 83:2481–2497.
- [123] Markou, M. and Singh, S. (2003b). Novelty detection: a review-part 2: Neural network based approaches. *Signal Processing*, 83:2499–2521.
- [124] Mishra, M. and Gupta, N. (2015). To detect outlier for categorical data streaming. *International Journal of Scientific & Engineering Research*, 6(5):1–5.
- [125] Moore, A., Lee, M. S., and Anderson, B. (1998). Cached sufficient statistics for efficient machine learning with large datasets. *Journal of Artificial Intelligence Research*, 8:67ï£i–91.
- [126] Moore, A. and Wong, W. K. (2003). Optimal reinsertion: A new search operator for accelerated and more accurate bayesian network structure learning. In *Proceedings of the 20^t h International Conference on Machine Learning*, pages 552–559.
- [127] Narita, K. and Kitagawa, H. (2008). Detecting outliers in categorical record databases based on attribute associations. In *Progress in WWW Research and Development*, pages 111–123, Berlin, Heidelberg. Springer Berlin Heidelberg.
- [128] Noto, K., Brodley, C., and Slonim, D. (2010). Anomaly detection using an ensemble of feature models. In *IEEE International Conference on Data Mining (ICDM)*, pages 953–958, Sydney, Australia.
- [129] Noto, K., Brodley, C., and Slonim, D. (2012). Frac: A feature-modeling approach for semi-supervised and unsupervised anomaly detection. *Data Mining Knowledge Discovery*, 25(1):109–133.
- [130] O'Reilly, C., Gluhak, A., Imran, M. A., and Rajasegarar, S. (2014). Anomaly detection in wireless sensor networks in a non-stationary environment. *IEEE Communications Surveys & Tutorials*, 16(3):1413–1432.
- [131] Otey, M. E., Ghoting, A., and Parthasarathy, S. (2006). Fast distributed outlier detection in mixed-attribute data sets. *Data Mining and Knowledge Discovery*, 12(2-3):203–228.
- [132] Otey, M. E., Parthasarathy, S., and Ghoting, A. (2005). An empirical comparison of outlier detection algorithms. In Proceedings of the International Workshop on Data Mining Methods for Anomaly Detection at ACM International Conference on Knowledge Discovery and Data Mining, SIGKDD, pages 1–8, Chicago, IL, USA.
- [133] Pang, G., Cao, L., and Chen3, L. (2016). Outlier detection in complex categorical data by modeling the feature value couplings. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, pages 1902 1908, New York, NY, USA.
- [134] Patcha, A. and Park, J.-M. (2007). An overview of anomaly detection techniques: Existing solutions and latest technological trends. *Computer Network*, 51(12):3448–3470.
- [135] Pawar, M. S., Amruta, D., and Tambe, S. N. (2014). A survey on outlier detection techniques for credit card fraud detection. *IOSR Journal of Computer Engineering*, 16(2):44–48.
- [136] Pawlak, Z. (1982). Rough sets. International journal of computer & information sciences, 11(5):341-356.

- [137] Phua, C., Alahakoon, D., , and Lee, V. (2004). Minority report in fraud detection: Classification of skewed data. *ACM SIGKDD Explorations Newsletters*, 6(1):50âĂŞ-59.
- [138] Phua, C., Lee, V. C. S., Smith-Miles, K., and Gayler, R. W. (2010). A comprehensive survey of data mining-based fraud detection research. http://arxiv.org/abs/1009.6119.pdf, pages 1–14.
- [139] Pimentel, M. A., Clifton, D. A., Clifton, L., and Tarassenko, L. (2014). A review of novelty detection. Signal Processing, 99:215–249.
- [140] Pradip, S. S., Robert, J. F., and Hamza, J. F. (2015). Information-theoretic outlier detection for large-scale categorical data. *International Journal of Computer Science and Mobile Computing*, 4(4):873–881.
- [141] Purankar, R. M. and Patil, P. (2015). A survey paper on an effective analytical approaches for detecting outlier in continuous time variant data stream. *International Journal of Engineering and Computer Science*, 4(11):14946–14949.
- [142] Ramaswamy, S., Rastogi, R., and Shim, K. (2000). Efficient algorithms for mining outliers from large data sets. In *Proceedings of the ACM International Conference on Management of Data, SIGMOD*, pages 427–438, Dallas, Texas, USA.
- [143] Ranshous, S., Shen, S., Koutra, D., Harenberg, S., Faloutsos, C., and Samatova, N. F. (2015). Anomaly detection in dynamic networks: A survey. *Wiley Interdisciplinary Reviews: Computational Statistics*, 7(3):223–247.
- [144] Rashidi, L., Hashemi, S., and Hamzeh, A. (2011). Anomaly detection in categorical datasets using bayesian networks. In *The Third International Conference on Artificial Intelligence and Computational Intelligence, Part II, AICI*, pages 610–619, Taiyuan, China.
- [145] Rassam, M. A., Maarof, M., and Zainal, A. (2012). A survey of intrusion detection schemes in wireless sensor networks. *American Journal of Applied Sciences*, 9(10):1636–1652.
- [146] Rassam, M. A., Zainal, A., and Maarof, M. A. (2013). Advancements of data anomaly detection research in wireless sensor networks: A survey and open issues. *Sensors*, 13(8):10087–10122.
- [147] Reddy, D. L. S., Babu, B. R., and Govardhan, A. (2013). Outlier analysis of categorical data using navf. *Informatica Economica*, 17(1):1–5.
- [148] Rezaei, A., Kasirun, Z. M., Rohani, V. A., and Khodadadi, T. (2013). Anomaly detection in online social networks using structure-based technique. In *Proceedings of the International Conference for Internet Technology and Secured Transactions, ICITST*, pages 619–622, London, UK.
- [149] Ritika, Kumar, T., and Kaur, A. (2013). Outlier detection in wsn: A survey. *International Journal of Advanced Research in Computer Science and Software Engineering*, 3(7):609 –617.
- [150] Rokhman, N., Subanar, and Winarko, E. (2016). Improving the performance of outlier detection methods for categorical data by using weighting function. *Journal of Theoretical and Applied Information Technology*, 83:327–336.
- [151] Rousseeuw, P. J. and Driessen, K. V. (1998). A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, 41:212–223.
- [152] Sagade, A. G. and Thakur, R. (2014). Excess entropy based outlier detection in categorical data set. *International Journal of Advanced Computational Engineering and Networking*, 2(8):56–61.
- [153] Said, A. M., Dominic, D. D., and Samir, B. B. (2013). Outlier detection scoring measurements based on frequent pattern technique. *Research Journal of Applied Sciences Engineering and Technology*, 6(8):1340–134.
- [154] Sari, A. (2015). A review of anomaly detection systems in cloud networks and survey of cloud security measures in cloud storage applications. *Journal of Information Security*, 6(02):142–154.
- [155] Sarma, D. S. and Sarma, S. S. (2015). A survey on different graph based anomaly detection techniques. *Indian Journal of Science and Technology*, 8(31):1–7.
- [156] Savage, D., Zhang, X., Yu, X., Chou, P., and Wang, Q. (2014). Anomaly detection in online social networks. *Social Networks*, 39:62–70.
- [157] Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J., and Williamson, R. C. (2001). Estimating the support of a high-dimensional distribution. *Neural Computing*, 13(7):1443–1471.
- [158] Seok, J. and Kang, Y. S. (2015). Mutual information between discrete variables with many categories using recursive adaptive partitioning. *Scientific reports*, 5:1–10.
- [159] Shahid, N., Naqvi, I. H., and Qaisar, S. B. (2015). Characteristics and classification of outlier detection techniques for wireless sensor networks in harsh environments: A survey. *Artificial Intelligence Review*, 43(2):193–228.
- [160] Shannon, C. E. (1948). A mathematical theory of communication. *The Bell Telephone System Technical publication*, 27(3):379–423.
- [161] Shukla, D. S., Pandey, A. C., and Kulhari, A. (2014). Outlier detection: A survey on techniques of WSNs involving event and error based outliers. In *Proceedings of the International Conference of Innovative Applications of Computational Intelligence on Power, Energy and Controls with their Impact on Humanity, CIPECH*, pages 113–116, Ghaziabad, India.
- [162] Shyu, M., Sarinnapakorn, K., Kuruppu-Appuhamilage, I., Chen, S., Chang, L. W., and Goldring, T. (2005). Handling nominal features in anomaly intrusion detection problems. In *Proceedings of the International Workshop on Research Issues in Data Engineering: Stream Data Mining and Applications*, pages 55–62, Tokyo, Japan.

- [163] Singh, K. and Upadhyaya, S. (2012). Outlier detection: applications and techniques. *International Journal of Computer Science Issues*, 9(1):307–323.
- [164] Smets, K. and Vreeken, J. (2011). The odd one out: Identifying and characterising anomalies. In *Proceedings of the SIAM International Conference on Data Mining, SDM*, pages 804–815, Arizona, USA.
- [165] Sodemann, A. A., Ross, M. P., and Borghetti, B. J. (2012). A review of anomaly detection in automated surveillance. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews,* 42(6):1257–1272.
- [166] Supriya, G. and Shinde, S. M. (2015). Outliers detection using subspace method: A survey. *International Journal of Computer Applications*, 112(16):20–22.
- [167] Suri, N. N. R. R., Murty, M. N., and Athithan, G. (2012). An algorithm for mining outliers in categorical data through ranking. In *The Proceedings of 12th international conference on hybrid intelligent systems (HIS), IEEE*, pages 247ï£i–252, Pune, India.
- [168] Suri, N. N. R. R., Murty, M. N., and Athithan, G. (2013). A rough clustering algorithm for mining outliers in categorical data. In *The Proceedings of 4th international conference on pattern recognition and machine intelligence (PReMI)*, pages 170ï£;–175, Kolkata, India.
- [169] Suri, N. N. R. R., Murty, M. N., and Athithan, G. (2014). A ranking-based algorithm for detection of outliers in categorical data. *International Journal of Hybrid Intelligent Systems*, 11:1–11.
- [170] Suri, N. N. R. R., Murty, M. N., and Athithan, G. (2016). Detecting outliers in categorical data through rough clustering. *Nat. Comput.*, 15:385ï£;–394.
- [171] Taha, A. and Hadi, A. S. (2013). A general approach for automating outliers identification in categorical data. In ACS/IEEE International Conference on Computer Systems and Applications (AICCSA), pages 1–8, Ifrane, Morocco.
- [172] Taha, A. and Hadi, A. S. (2016). Pair-wise association for categorical and mixed attributes. *Information Sciences*, 346:73–89.
- [173] Taha, A. and Hegazy, O. (2010). A proposed outliers identification algorithm for categorical data sets. In *Proceedings* of *International Conference on Informatics and Systems, INFOS*, pages 1–5, Cairo, Egypt.
- [174] Wang, Y. (2008). Statistical Techniques for Network Security: Modern Statistically-Based Intrusion Detection and Protection. IGI Global, New York, NY, USA.
- [175] Wang, Y. and Xu, W. (2018). Leveraging deep learning with lda-based text analytics to detect automobile insurance fraud. *Decision Support Systems*, 105:87–95.
- [176] Wei, L., Qian, W., Zhou, A., Jin, W., and Yu, J. X. (2003). Hypergraph-based outlier test for categorical data. In *Proceedings of the ACM International Conference on Knowledge Discovery and data Mining, SIGKDD*, pages 399–410, Washington, DC, USA.
- [177] Weller-Fahy, D. J., Borghetti, B. J., and Sodemann, A. A. (2015). A survey of distance and similarity measures used within network intrusion anomaly detection. *IEEE Communications Surveys & Tutorials*, 17(1):70–91.
- [178] West, J. and Bhattacharya, M. (2016). Intelligent financial fraud detection: A comprehensive review. *Computers & Security*, 57:47–66.
- [179] Wu, S. and Wang, S. (2011). Parameter-free anomaly detection for categorical data. Machine Learning and Data Mining in Pattern Recognition. Lecture Notes in Computer Science, 6871:112âÅŞ-126.
- [180] Wu, S. and Wang, S. (2013). Information-theoretic outlier detection for large-scale categorical data. *IEEE Transactions on Knowledge Data Engineering*, 25(3):589–602.
- [181] Yassin, W., Udzir, N. I., Muda, Z., and Sulaiman, N. (2013). Anomaly-based intrusion detection through k-means clustering and naives Bayes classification. In *Proceedings of the International Conference on Computing and Informatics, ICOCI*, pages 298–303, Bandung, Indonesia.
- [182] Yu, J. X., Qian, W., Lu, H., and Zhou, A. (2006). Finding centric local outliers in categorical/numerical spaces. *Knowledge and Information Systems*, 9:309–338.
- [183] Yu, R., Qiu, H., Wen, Z., Lin, C.-Y., and Liu, Y. (2016). A survey on social media anomaly detection. http://arxiv.org/pdf/1601.01102.pdf, pages 1–24.
- [184] Zhang, J. (2013). Advancements of outlier detection: A survey. ICST Transactions on Scalable Information Systems, 13(1):1–26.
- [185] Zhang, Y., Meratnia, N., and Havinga, P. (2010). Outlier detection techniques for wireless sensor networks: A survey. *IEEE Communications Surveys & Tutorials*, 12(2):159–170.
- [186] Zhao, X., Liang, J., and Cao, F. (2014). A simple and effective outlier detection algorithm for categorical data. *International Journal of Machine Learning and Cybernetics*, 5:469–477.
- [187] Zijlstra, W. P., van der Ark, L. A., and Sijtsma, K. (2011). Outliers in questionnaire data: Can they be detected and should they be removed. *Journal of Educational and Behavioral Statistics*, 36:186–212.
- [188] Zimek, A., Schubert, E., and Kriegel, H.-P. (2012). A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining*, 5(5):363–387.