

Outlier Detection Based on Surfeit Entropy for Large Scale Categorical Data Set

Neha L. Bagal

PVPIT, Bavdhan, Pune, Savitribai Phule Pune University, India

Abstract: *Number of methods based on classification, clustering, frequent patterns and statistics has been proposed to collect meaningful information by removing unwanted data. Information theory uses statistical approach to achieve its goal. The outlier detection from unsupervised data sets is more difficult task since there is no inherent measurement of distance between objects. Here in this work, we proposed a novel framework based on information theoretic measures for outlier detection in unsupervised data with the help of Max/Surfeit Entropy. In which we are using different information theoretic measures such as entropy and dual correlation. Using this model we proposed SEB-SP outlier detection algorithm which do not require any user defined parameter except input data. We have also used the concept of weighted entropy. Our method detects outliers better than existing approach.*

Keywords: Outlier detection, surfeit entropy, weighted entropy, dual correlation

1. Introduction

Outlier detection is an important research problem which is researched within various research areas and application domains. Outlier detection method detects distinct, exceptional and inconsistent objects with respect to the majority data in a given input data sets. Number of outlier detection techniques has been specifically developed for some application areas [1]. While others are many times known as anomaly detection in different literature. Outlier detection can also utilize in scientific research work for analyzing data and knowledge discovery in astronomy, oceanography, chemistry, biology and other applications.

Outlier arises due to faults in mechanical systems, system's behavioral changes, mankind errors, fraudulent nature, and instrument errors [1]. Detection of these outliers helps to identify system faults before they affect intensively with outcomes of a system. The techniques or algorithms used in outlier detection methods are varied notably which are mainly dependent upon the characteristics of data sets to be worked with. If we classify the existing methods for outlier detection according to the availability of labels in the training data sets, there are three broad categories [1]:

1. Supervised,
2. Semi-supervised, and
3. Unsupervised approaches.

In supervised or the semi-supervised approaches all need to be trained before use, while the unsupervised approach do not include the training phase. For a supervised approach a training set should be provided with labels for anomalies as well as labels of normal objects, in contrast with the training set with normal object labels alone required by the semi-supervised approach. On the other side, the unsupervised approach does not require any object label information.

Above three approaches have different prerequisites, limitations and uses different data sets with different amounts of label information. These approaches are discussed [8] in detail below.

Supervised outlier detection approach uses labeled objects belonging to the normal and outlier classes to learn the classifier and assign appropriate labels to test objects.

Semi-supervised outlier detection approach firstly learns normal behavior from given training data set of normal objects and then calculates the similarity of test objects.

Unsupervised outlier detection approach detects outliers in unlabeled data set. Considering that the most of the objects in data set are normal. This approach is applied to various kinds of outlier detection methods and data sets. We mainly concentrate on unsupervised approach in this paper. To use supervised and semi-supervised approach one must first label the training data sets. When we consider large data sets or high dimensional data then labeling will be tedious and time consuming task.

This paper is organized as follows: section II describes the objective and associated challenges. Section III consists of literature survey of existing work with comparisons. Some background details and formulation is given in section III. Section IV describes proposed system with mathematical models. Section V gives conclusion and Section VI describes references.

2. Challenges and Objectives

- 1) Outliers are the patterns that do not conform to expected normal behavior. Therefore we define normal behavior of objects and declare any observation which does not belong to normal behavior as an anomaly but several factors make this simple approach very challenging.
- 2) Defining a boundary between normal and anomalous behavior is not simple task, Adaption in anomalous observations to appear like normal, difference in exact notion of outliers for different domains. This makes outlier detection problem more complex and difficult.
- 3) Unsupervised method are applicable only on numerical data sets, however they cannot be used to deal with categorical data.
- 4) Using formal definition of outlier our aim is to develop effective and efficient method that can be used to detect

outliers in large scale categorical unsupervised data sets for real applications.

- 5) We have combined entropy and dual correlation with attribute weighting resulting into weighted surfeit entropy where entropy computes uncertainty and dual correlation measures mutual information or attribute relation.
- 6) Different user defined intrinsic and decision parameters are required for outlier detection problem. Thus results are based on correct values of these parameters.

3. Literature Survey

Methods designed for unsupervised outlier detection in categorical data can be grouped into four categories as follows.

A. Proximity Based Methods

In order to understand this concept, it is the method which measures compactness of objects in terms of distance / density. ORCA [10] and CNB [11] are different algorithms for outlier detection in categorical data. ORCA uses hamming distance and CNB uses common neighbor set. These two methods are not useful for high dimensional data because of difficulty in choosing the distance or density as well as high time and space complexity.

B. Ruled Based Methods

Rule-based methods use the concept of frequent items from association-rule mining. This method considers the frequent or infrequent items as a data set. Objects with few frequent items or many infrequent items are more likely to be considered as anomalous objects than others.

Frequent pattern outlier factor [12] and Otey's algorithm [13] are two well known ruled based techniques. The FIB algorithm includes an initial computation of the set of frequent patterns, using a predefined minimum support rate. All support rates of associated frequent patterns are summed up for each object as the outlier factor of this object. While Otey's algorithm, begins with computation of infrequent items from data set. Outlier factor is calculated using the same. Objects with largest scores are treated as outliers.

C. Other Methods

Random walk, Hyper-graph theory [8] methods are implemented using several approaches. In random walk method [14], outliers are the object who has the low probability to combine with neighbor. That means object remain in their state. In method [15] relationships are considered and mutual dependence based local outlier factor is proposed to detect outliers. There are many other methods cluster based local outlier detection method, classification based method.

In literature several methods have been proposed for outlier detection using information theoretic measures.

- 1) An information theoretic method is proposed for anomaly detection in audit data sets in [2] using measures like entropy, conditional entropy, relative entropy & information gain. This method aims to identify

outliers in the univariate audit data set, where regularity is characterized but not the attribute relation.

- 2) Information theoretic outlier detection in large scale categorical data holoentropy –sum of entropy and total correlation. This method gives optimal solution to outlier detection by using ITB-SP algorithm.

ITB-SP Method computes holoentropy as follows:

$$HLx(y) = Hx(y) + Cx(y) \tag{1}$$

After computing holoentropy, weighted holoentropy and differential entropy is calculated using (2) and (3) to get anomaly set AS from data set.

$$Wx(y) = \sum_{i=1}^m Wx(yi)Hx(yi) \tag{2}$$

$$Hx(xo) = Wx(y) - Wx\{xo\}(y) \tag{3}$$

Table 1: Comparison of systems

Parameter	CNB	ORCA	FPOF	ITB-SP
Approach	Proximity based	Proximity based	Ruled based	Information theoretic based
Method	Distance	Distance	Item set frequency	High Dimensional categorical data
Input Data Set	Low dimensional categorical data	High dimensional data in random	Low dimensional Numeric data	High dimensional categorical data
Required parameters	M, sim, k	K, M	Minfreq, maxlen, M	Number of outliers o
Output Data Set	outliers	O-outliers	Value of FPOF, FP-outliers	OS-outlier set
Complexity	$O(n^2(k+S(\theta) + q) + n(k + M))$	$O(n^2q)$	$O(n(2^{T-1}))$	$=O(nm)$

Table 1, shows the comparison between different outlier detection methods using parameters like approach, type, input data set, output data set, complexity and user defined parameters.

D. Necessity of modified outlier detection:

Most of the existing systems are depends on user defined parameters and very few methods are dealing with unsupervised categorical data. Therefore a method should exist which will be able to deal large scale categorical data without requirement of any user defined parameter. There is requirement of method which will perform outlier detection using joint correlation between attributes.

4. Problem Formulation

In this section we first look at how entropy and dual total correlation can be used to capture similarity between outlier candidates. We are proposing Weighted Surfeit Entropy and formulate the outlier detection problem.

A. Entropy: Entropy is measure of information and uncertainty of a random variable.

Let X be the set of n objects $\{x_1, x_2, x_3, \dots, x_n\}$, each x_i for $1 \leq i \leq n$ being a vector of categorical attributes $[y_1, Y_2, y_3, \dots, y_m]^T$ where m number of its attributes. Now based on chain rule of entropy [4], Entropy of y denoted as $Hx(y)$ can be written as follows.

$$Hx(y) = Hx(y_1, y_2, \dots, y_m)$$

$$= \sum_{i=1}^m Hx(y_i | y_1, \dots, y_{i-1})$$

$$= Hx(y_1) + Hx(y_2 | y_1) + \dots + Hx(y_m | y_1, \dots, y_{m-1})$$

(4)

Where

$$Hx(y_m | y_1, \dots, y_{m-1}) = - \sum_{y_m, y_{m-1}, \dots, y_1} p(y_m, y_{m-1}, \dots, y_1) \log p(y_m | y_1, \dots, y_{m-1}).$$

Entropy of dataset decreases significantly with removal good outlier candidates.

B. Total Correlation:

It is defined as summation of mutual information of multivariate discrete random vector y , [8,1] and it is denoted as $C_x(y)$. Total correlation is based on Watanabe's proof. Total correlation can be expressed as :

$$C_x(y) = \sum_{i=1}^m Hx(y_i) - Hx(y)$$

(5)

C. Dual total correlation:

The dual total correlation [17] calculates the amount of entropy present in Y beyond the sum of the entropies for each variable conditioned upon all other variables. The dual total correlation is also called as the surfeit entropy and the binding information. In this paper we describe dual total correlation as $Ex(Y)$ and expressed as

$$Ex(y) \equiv (\sum_{x_i \in X} Hx \setminus x_i(y)) - (n - 1)Hx(y)$$

(6)

Where n is number of attributes.

To weight the entropy of each attribute, we are using a reverse function of the entropy, as follows:

$$Wx(y_i) = 2 \left(1 - \frac{1}{1 + \exp \left(\frac{1}{Hx(y_i)} \right)} \right)$$

(7)

The weighted Surfeit entropy is defined as follows:

Definition 1: The weighted surfeit entropy $EW_x(Y)$ is the sum of weighted entropy on each attribute of the random vector Y .

$$EW_x(Y) = \sum_{i=1}^m Wx(y_i) Hx(y_i)$$

(8)

Outliers are detected by minimizing the surfeit entropy through the removal of outlier candidates; Proposed strategy Have weighting the entropy of each individual attribute in

order to give more importance to those attributes with small entropy values.

D. Formal definition of outlier detection:

We are using weighted surfeit entropy for outlier detection outliers .We consider that set of outlier candidates is the best if entropy of dataset significantly decreases with its removal from dataset.

Definition 2: X be a given dataset with n objects and a subset $Out(o)$ is defined as the set of outliers if it minimizes the weighted surfeit entropy of dataset X with o objects removed.

E. Differential surfeit entropy:

Definition 3: Given an object x_o of X , the difference of weighted surfeit entropy $e_x(x_o)$ between the data set X and the data set $X \setminus \{x_o\}$ is defined as the differential surfeit entropy of the object x_o .

$$e_x(x_o) = Wx(y) - Wx \setminus \{x_o\}(y)$$

(9)

F. Outlier factor:

Outlier factor is a measure of how likely x_o is an outlier. An object x_o with a large outlier factor value is more likely to be an outlier than an object with a small value. Outlier factor of an object x_o is denoted as $OF(x_o)$ is defined as :

$$OF(x_o) = \sum_{i=1}^m OF(x_o, i)$$

(10)

5. Proposed Approach

Our proposed approach is based on weighted entropy and differential entropy which can be calculated using equation (8) and (9). System will take data set file of format .CSV and gives output file with outliers removed.

System Architecture

To address the problem discussed above in need of effective outlier detection in unsupervised data set. A proposed methodology with surfeit entropy and to deal with large scale categorical data is considered as shown in Fig. 1

1) GUI Handler:

It provides following functionality:

- File selector (CSV File)
- Display for Attributes
- Display for Outliers (Outcome)

2) File Processor:

It will handle following tasks:

- Separate objects and attributes.
- Saving outlier results.

3) Outlier Detector:

It will handle following tasks:

- Calculate Entropy
- Calculate Dual total Correlation
- Calculate weighted surfeit entropy
- Calculate Outlier factor
- Getting outlier set
- Getting data set file with removal of attributes

4) Report generator:

- Generate Report
- Generate comparison model using graphs

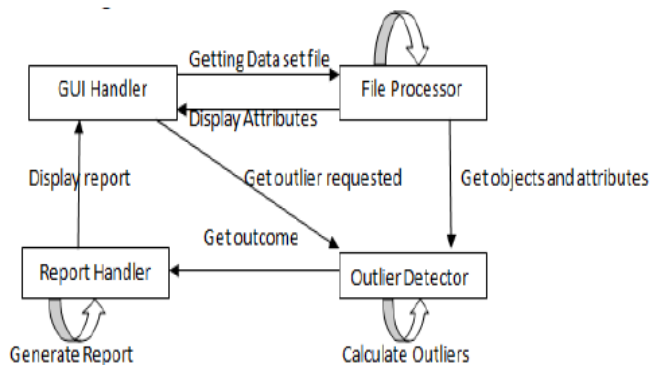


Figure 1: System Architecture

Mathematical model:

The proposed concept is constructed on the assumption that elimination of outliers will improve the purity of data set and reduces $EWx(y)$. When a normal object is removed from the data set, the value of $EWx(y)$ should increase. Thus, the objects with positive $e(x_i)$ are defined as the outlier candidate set (OS). The objects with non-positive $e(x_i)$ are defined as elements of normal object set (NS).

$$OS = \{Xi | e(Xi) > 0\} \text{ And } NS = \{Xi | e(Xi) \leq 0\}$$

SEB-SP Algorithm for outlier detection:

In this section, we have derived Surfeit entropy based single pass greedy algorithm for outlier detection. In this algorithm outlier factors are computed only once, and the o objects with largest values are identified as outliers. This algorithm is parameter-less as we do not need to provide any user defined parameters.

Algorithm: SEB-Single Pass

1. **Input:** data set X
2. **Output:** Outlier set S
3. Compute $w(y_0)$ for $(1 \leq i \leq m)$ by (7)
4. Set OS=null
5. **for** $i=1$ to n **do**
6. Compute OF(xi) and obtain OS by (10)
7. **end for**
8. Build S by searching in OS

Greedy approach is used to find out outliers from input data set in above algorithm. Algorithm firstly computes weighted entropy for each attribute. Then entropy of each attribute is updated. The attribute entropy is always changes when outliers are detected and removed from the data set, then calculates outlier factor for each attribute and get the largest OF set which will convert to OS (Outlier set). After that set S will be built. Complexity of the algorithm is $O(nm)$, as we are not using any searching algorithm.

6. Conclusion

This paper discusses many outlier detection methods based on information theory. We are proposing novel method which will overcome limitations of previous methods. This

paper formulates outlier detection as an optimization problem and proposed a practical, unsupervised, parameter less algorithm for detecting outliers in large-scale categorical data sets. Effectiveness of our approach results from a new concept of surfeit entropy. The efficiency of our algorithms results from the outlier factor function derived from the differential entropy. In particular, in our proposed approach dual total correlation and surfeit entropy works more effectively to remove outlier from large scale categorical attributes.

References

- [1] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly Detection: A Survey," ACM Computing Surveys, vol. 41, no. 3, pp. 1-58, 2009.
- [2] V.J. Hodge and J. Austin, "A Survey of Outlier Detection Methodologies," Artificial Intelligence Rev., vol. 22, no. 2, pp. 85- 126, 2004.
- [3] E.M. Knorr and R.T. Ng, "Algorithms for Mining Distance-Based Outliers in Large Data Sets," Proc. 24rd Int'l Conf. Very Large Data Bases (VLDB '98), 1998.
- [4] S.R. Gaddam, V.V. Phoha, and K.S. Balagani, "K-Means+ID3: A Novel Method for Supervised Anomaly Detection by Cascading K-Means Clustering and ID3 Decision Tree Learning Methods," IEEE Trans. Knowledge and Data Eng., vol. 19, no. 3, pp. 345-354, Mar. 2007.
- [5] V. Chandola, Banerjee, and Kumar, "Anomaly Detection: A Survey", ACM Computing Surveys, vol. 41, no. 3, pp. 1-58, 2009.
- [6] V. Hodge and J. Austin "A Survey of Outlier Detection Methodologies," Artificial Intelligence Rev, vol. 22, no. 2, pp. 85-126, 2012.
- [7] J. Zhang," Advancements of Outlier Detection: A Survey," ICST Transactions on Scalable Information Systems, Vol.3, no. 1, March 2013.
- [8] Sheu Wu, Member IEEE, and Shengrui Wang, Member IEEE, "Information-Theoretic Outlier Detection for Large-Scale Categorical Data," IEEE transactions on knowledge and data engineering, vol. 15, no. 3, march 2013.
- [9] T. Cover and J. Thomas, "Elements of Information Theory" John Wiley & Sons, 1991, pp.12-21.
- [10] Ayman Taha, Ali Hadi," A General Approach for Automating Outliers Identification in Categorical Data, IEEE conference on Computer System & applications, pp.1 – 8, May 2013.
- [11] S. Li, R. Lee, "Mining Distance-Based Outliers from Categorical Data," Proc. IEEE Seventh International Conference, Data Mining Workshops (ICDM '07), 2007.
- [12] S.D. Bay "Mining Distance-Based Outliers in Near Linear Time with Randomization and a Simple Pruning Rule," Proc. Ninth ACM SIGKDD International Conf. Knowledge Discovery and Data Mining (KDD '03), 2008.
- [13] Z. He, X. Xu, Z.J. Huang, "FP-Outlier: Frequent Pattern Based Outlier Detection," Computer Science and Information Systems, vol. 2, pp. 103-118, 2005.
- [14] M.E. Otey, A. Ghoting, "Fast Distributed Outlier Detection in Mixed-Attribute Data Sets," Data Mining and Knowledge Discovery, vol. 12, pp. 203-228, 2006.

- [15] H.D.K. Moonesignhe, Tan, "Outlier Detection Using Random Walks," Proc. IEEE 18th Int'l Conf. Tools with Artificial Intelligence (ICTAI '09), 2009.
- [16] W. Qian, H. Lu, and A. Zhou, "Finding Centric Local Outliers in Categorical/Numerical Spaces, Knowledge and Information Systems, vol. 8, no. 3, pp 309-338, 2006.
- [17] Nicholas Timme, Wesley Alford, Benjamin Flecker, and John M. Beggs, "Multivariate information measures: an experimentalist's perspective", 28 Nov 2011.
- [18] S. Watanabe, "Information Theoretical Analysis of Multivariate Correlation," IBM J. Research and Development, vol. 4, pp. 66-82, 1960.
- [19] <http://www.data setgenerator.com/>, 2011.