

An Evaluation of Detection of Outliers by Reverse Nearest Neighbors Method

A.Prashanth¹, K.A.M.Sushma², P.Srinivas Rao³

¹M.Tech ,CSE, Jayamukhi Institute Of Technological Sciences, Warangal, India

²Assistant professor, CSE, Jayamukhi Institute Of Technological Sciences, Warangal, India

³Associate professor, CSE, Jayamukhi Institute Of Technological Sciences, Warangal, India

ABSTRACT: Outlier Detection in high dimensional information goes into a rising system in today's inspection in the region of data mining. Data that is different or errant from normal data set are recognized by outlier detection. Unusual data records because of some data errors can be treated as outliers typically detecting outliers and investigating large data sets recognizes the problem such as medical problems, a structural defect, and investigational errors. This paper focuses the different methods for detection of anomalies. In order to handle the difficulties related to outlier detection because of uncertain data, outlier detection technique based on the AntiHub term is used.

KEYWORDS- Data stream, Data mining, outlier detection.

I. INTRODUCTION

In spite of the huge measure of information being gathered in numerous exploratory and business applications, specific occasions of hobbies are still very uncommon. These uncommon occasions, regularly called exceptions or irregularities, are characterized as occasions that happen occasionally (their recurrence ranges from 5% to under 0.01% relying upon the application). Discovery of exceptions (uncommon occasions) has as of late picked up a great deal of consideration in numerous areas, extending from video observation and interruption identification to fake exchanges and coordinate advertising. For instance, in video observation applications, video directions that speak to suspicious and/or unlawful exercises (e.g. recognizable proof of movement violators out and about, discovery of suspicious exercises in the region of articles) speak to just a little divide of all video directions. Thus, in the system interruption discovery area, the quantity of digital assaults on the system is

regularly a little portion of the aggregate system movement. In spite of the fact that exceptions (uncommon occasions) are by definition rare, in each of these illustrations, their significance is entirely high contrasted with different occasions, making their identification critical. Information excavating strategies produced for this issue depend on both managed and unsupervised learning. Regulated learning routines commonly assemble an expectation model for uncommon occasions in light of named information (the preparation set), and utilize it to arrange every occasion [7-8]. The significant disadvantages of regulated information mining strategies include: (1) need to have marked information, which can be to a great degree tedious for genuine applications, and (2) powerlessness to identify new sorts of uncommon occasions. Interestingly, unsupervised learning systems commonly don't require marked information and distinguish exceptions as information focuses that are altogether different from the typical (greater part) information in light of some measure [9].

These strategies are ordinarily called exception/irregularity recognition procedures, and their prosperity relies upon the decision of closeness measures, highlight choice and weighting, and so on. They have the upside of distinguishing new sorts of uncommon occasions as deviations from typical conduct, yet then again they experience the ill effects of a conceivable high rate of false positives, basically since already concealed (yet ordinary) information can be likewise perceived as exceptions/oddities. Regularly, information in numerous uncommon occasions applications (e.g. system movement observing, video observation, web use logs) arrives persistently at a tremendous pace in this way representing a noteworthy test to break down it [9]. In such cases, it is imperative to settle on choices rapidly and precisely. On the off chance that

there is a sudden or startling change in the current conduct, it is fundamental to distinguish this change as quickly as time permits. Expect, for instance, there is a PC in the neighborhood that uses just set number of administrations (e.g., Web activity, telnet, ftp) through comparing ports. Every one of these administrations relate to specific sorts of conduct in system activity information. On the off chance that the PC all of a sudden begins to use another administration (e.g., ssh), this will positively resemble another sort of conduct in system activity information. Henceforth, it will be attractive to identify such conduct when it shows up particularly since it might frequently relate to unlawful or nosy occasions. Indeed, even for the situation when this particular change in conduct is a bit much nosy or suspicious, it is imperative for a security examiner to comprehend the system activity and to redesign the idea of the typical conduct. Further, on-line recognition of irregular conduct and occasions additionally assumes a huge part.

II. RELATED WORKS

In recent time it is observed that the distribution of points' reverse-neighbor counts becomes skewed in high dimensions, which results in the phenomenon of hubness [1]. Authors also discussed that the how antihub appear very infrequently in k-NN lists of other points. They also discussed the connection between the antihubs and existing unsupervised outlier detection [1].

Here provided the role of reverse nearest neighbor counts in problems concerning unsupervised outlier detection. The main focus is given on the unsupervised outlier-detection methods and the hubness phenomenon in high dimensionality. Extended the work of antihubs to the large values of k and explored the relation between the hubness and data sparsity based on the unsupervised outlier detection. The extension of antihubs improves the discrimination in the outlier scores. The existence of hubs and antihubs in high-dimensional data is relevant to machine-learning techniques from various families: supervised, semi-supervised, as well as unsupervised. Here only unsupervised method is used, it does not give accurate result as compared to the other methods.

H.-P. Kriegel, M. Schubert, and A. Zimek [4] has proposed angle based outlier detection (ABOD). Outlier detection in high-dimensional data uses the variances of a measure over angles between the different vectors of data objects. In ABOD technique, used the properties of the variances to actually take advantage of high dimensionality and found to be less sensitive to the increasing dimensionality of a dataset. This technique is less efficient than the classic distance based methods. The disadvantage is only angle based is used not the classic distance-based methods.

The LOF compare the local density of instances with the densities of its neighborhood instances. After that it assign the outlier scores to given data objects. If LOF score equal to ratio of average local density of k nearest neighbor of instance and local density of data instance itself then data instance is considered to be normal and not as an outlier. Local density of instances is computed by finding radius of small hyper sphere centered at the data instance after that dividing volume of k [5], i.e. k nearest neighbor and volume of hyper sphere. In this assign a degree to each object to being an outlier known as local outlier factor [5].

Objects are isolated depending on the surrounding neighborhood, instances lying in dense region are normal objects [5], if their local density is similar to their neighbors and objects are outlier if their local density lower than its nearest neighbor [5]. It is a critical or lengthy process as compared to the distance based methods. The antihub2 method is unsupervised outlier detection method used for anomaly detection in high dimensional dataset. Anomaly detection in high dimensional data exhibits that as dimensionality increases there exists hubs and antihubs [6]. Hubs are the point that frequently occurs in k nearest neighbors. Antihubs are the point that occurs infrequently in nearest neighbors list. In this paper authors have refined the antihub method to refine the outlier scores of a point produced by the antihub method by considering the nk scores of the neighbors of the data point. Discrimination of outlier scores produced by Antihub2 acquires longer period of time with larger number of iterations [6]. Because of this recursive AntiHub2 method was introduced to improve the computational complexity

of discriminating the outlier scores using less number of iterations to detect accurate outliers in high dimensional data[6].

III. PROPOSED APPROACH

An outlier detector identifies statistics items that don't confirm an anticipated sample or different records items in facts stream. The detection of outlier enables in discovery of surprising information in records circulate. System version for outlier detection technique is proven in fig.1. It suggests fundamental phases for the procedure of identity of outliers and its additives. These are explained as follows:

A. Database series and pre-processing

For the input to the machine, datasets are collected from the UCI depository sets. In the device module widespread databases are used and the databases are supervised databases. Supervised dataset incorporates class labels in line with the data type. Class labels are assigned at the basis of different attributes of the dataset. That manner datasets already characterized into different lessons based at the data type. As in keeping with the form of dataset elegance labels according to magnificence kind are present. As part of preprocessing, the missing values inside the databases are filled with the value zero or as null.

The first pre-processing technique used is data cleansing. Data cleaning is executed to discover the lacking values in the data file. Also data cleansing is carried out to find out inconsistent information. The dataset used for outlier detection contained missing values and inconsistent information, for these statistics cleaning is achieved. As a few attribute does now not contain any cost, for them as part of statistics cleansing, values are inserted like 0 and null.

The every other part of pre-processing is facts transformation. It consists of conversion of data values aptitude in the shape of dataset into the data layout for destination device requiring records from source device consisting facts. Data transformation carries information mapping, which maps information from supply to vacation spot system. As the data is to be had on UCI depository, the data set is

mapped into the hotspot machine using data transformation technique. Hence, the datasets are equipped for the following step.

B. AntiHub1 Method

The approach AntiHub1 is primarily based on the ODIN method. ODIN method makes use of normal scores of outliers through analyzing the closest neighbor be counted for the unique factor. The ordered facts set D is given as an enter. Number of acquaintances and distance from unique point is supplied as input. Temporary variable AntiHub rankings are used for comparing cutting-edge discrimination score and raw outlier score. For each enter 'Sn' is computed w.r.t. Dist and statistics set (d/n) .

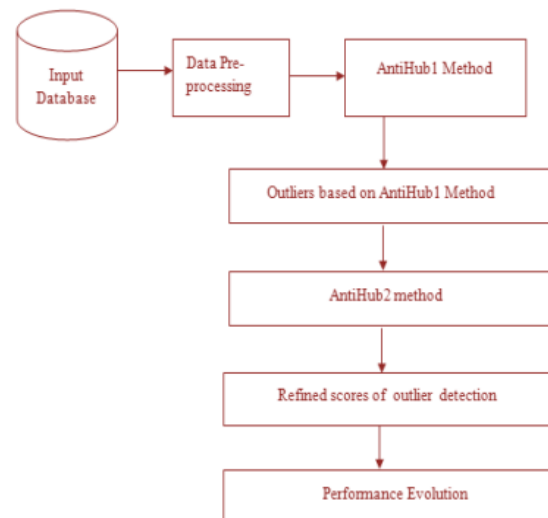


Fig. 1 Flowchart for outlier detection

The steps for AntiHub1 algorithm [13] are as follows:

1. For each point 'i' belongs to 1 to n,
2. $N_x(x) = D \setminus x_i$
3. $s = f(t)$

The steps for algorithm are described:

- Distance measured in the form of variable dist and ordered dataset as D is given as input for each point.
- For each point, distance from nearest neighbour is computed as $N_x(x)$ with respect to data set and dist variable.

- The score of outlier for point 'x' is computed using monotone function and stored in the vector 's' as an output.

The Algorithm 1 AntiHub produces output in the form of vector as the outlier score of point x from data set D. Outliers are collected as output in the form of percentage outlier values.

C. AntiHub2 Method

As normalization is not applied on the outlier score, the outliers collected are not refined. To tackle weaknesses of AntiHub1, a simple heuristic method AntiHub2 is applied to the outlier score produced by the AntiHub1 method. AntiHub2 refines outlier scores produced by the AntiHub1 method by considering the N_k scores of the neighbors of x.

For improvement in discrimination of scores that AntiHub2 introduces compared to AntiHub, for each point x, AntiHub2 proportionally adds the sum of N_k scores of the k nearest neighbors of x.

Input for AntiHub2 method contains measured distance for each point, ordered data set, number of neighbours and ratio of outliers for maximizing discrimination. Temporary variables are used for getting outlier score which are current discrimination score, current raw outlier scores, antihub score and sums of nearest neighbours' scores. AntiHub2 method is implemented for the refinement of the scores of outliers produced by AntiHub1 method.

The steps for AntiHub2 algorithm [13] are as follows:

1. For each point 'i' belongs to 1 to n,
2. AntiHub score $a = \text{AntiHub}(D, k)$
3. $anni =$ summation of indices of k nearest neighbour
4. $disc = 0$
5. For each 'i' from 1 to n
6. $ct = ai + ann$
7. $cdisc = \text{discscore}(ct, p)$
8. If $cdisc > disc$
9. $t = ct$
10. $s = f(t)$

The steps for algorithm are described:

- Distance measured in the form of variable dist, ordered dataset as D and number of neighbours is given as input for each point.
- Ratio of outliers to maximize discrimination and search parameter for each step is initialized applicable to every point.
- Array for ordered dataset and number of neighbours is initialised as AntiHub(Dike).
- For each point sum of nearest neighbours' AntiHub scores as temporary variable 'ann' is calculated and stored. For each step from 0 to 1, loop is carried out, for each point.
- Raw outlier score 'ct' is calculated using proportion and point sum of nearest neighbours' AntiHub scores.
- Then the value raw outlier score 'ct' and ratio of outliers to maximize discrimination i.e. 'p' is transferred to temporary variable 'disc'.
- Comparison for temporary variables 'cdisc' and 'disc' is carried out, and if $cdisc > disc$ then current raw outlier score 't' is set as raw outlier score 'ct'.

The score of outlier for point 'x' is computed using monotone function and stored in the vector 's' as an output. The second method considers the scores of neighbours for point x. And then adds sum of scores of nearest neighbour. To find aggregate of neighbours' scores summation is calculated. The discrimination scores compared using discScore parameter is provided to output vector as an outlier score.

IV. CONCLUSION

In this paper, comparative influences of basic parameters dependences on outlier detection are deliberated. Influence for values and datasets and their inter-relationship are also recognized. To conclude, this paper aids to comprehend the fact about complete analysis of nature of task is to be modelled prior to the algorithmic choice for outlier detection.

REFERENCES

- [1] Milos Radovanovi, Alexandros Nanopoulos and Mirjana Ivanovi, "Reverse Nearest Neighbors

inUnsupervised Distance-Based Outlier Detection”, IEEETransactions On knowledge And Data Engineering.Transactions, Vol. 27, No. 5, May 2015.

[2] Edwin, Raymond, “Distance based outliers: algorithmsand applications”, Springer- verlag, 2008.

[3] AlexandrosNanopoulos, YannisTheodoridis, YannisManolopoulos, “C2P: Clustering based on Closest Pairs”,Proceedings of the 27th VLDB Conference, Roma, Italy,2011.

[4] H.-P. Kriegel, M. Schubert, and A. Zimek, “Angle-basedoutlier detection in high-dimensional data, ” in Proc 14thACM SIGKDD Int. Conf. Knowl. Discovery DataMining, 2008, pp. 444–452.

[5] K. Zhang, M. Hutter, and H. Jin, “A new local distancebased outlier detection approach for scattered real-worlddata, ” in Proc 13th Pacific-Asia Conf on KnowledgeDiscovery and Data Mining (PAKDD), pp. 813–822.2009.

[6] J.Michael Antony Sylvia, Dr. T. C. Rajakumar Recursiveantihub “outlier Detection in High Dimensional Data.”Vol-2, Issue-8 PP. 1269-1274 global journal of research,2015.012.

[7] W. Lee, S. Stolfo,“Data mining approaches for intrusiondetection”, Proc. of the 7th USENIX security symposium,1998.

[8] E. Bloedorn, et al.,“Data Mining for Network IntrusionDetection: How to Get Started”, MITRE Technical Report,August 2001.

[9] A.K. Jones, R.S. Sielken,“Computer System IntrusionDetection: A Survey. Technical report”, University of VirginiaComputer Science Department, 1999.

[10] M. Masud, Q. Chen, L. Khan,J. Gao and J. Han “Classification andAdaptive Novel Class Detection of Feature Evolving Data Streams”,IEEE Trans. Knowl. Data Eng., vol. 25, no. 7, July 2013.

[11] S. Ahmed Shaikh and H. Kitagawa “Continuous Outlier Detection onUncertain Data Streams”, IEEE Ninth International Conference onIntelligent Sensors, Sensor Networks and Information

Processing(ISSNIP) Symposium on Information Processing Singapore, 21–24April 2014 .

[12] Bo Liu, Yanshan Xiao, Philip S. Yu, ZhifengHao, and LongbingCao, “An Efficient Approach for Outlier Detection with ImperfectData Labels”, IEEE Trans. Knowl. Data Eng., vol. 26, no. 7, July2014.

[13] Milos Radovanovic, AlexandrosNanopoulos, and MirjanaIvanovi“Reverse Nearest Neighbors in Unsupervised Distance-BasedOutlier Detection”, IEEE Transactions On Knowledge And DataEngineering, Vol. 27, No. 5, May 2015.

[14] UCI Machine Learning Repository [Online]. Available:<http://archive.ics.uci.edu/ml/datasets.html>.

[15] Bo Liu, Yanshan Xiao, Philip S. Yu, ZhifengHao, and LongbingCao, “An Efficient Approach for Outlier Detection with ImperfectData Labels”, IEEE Trans. Knowl. Data Eng., vol. 26, no. 7, July2014.