







Outlier Detection in Multivariate Time Series Data Using a Fusion of K-Medoid, Standardized Euclidean Distance and Z-Score

Nwodo Benita Chikodili^(✉) , Mohammed D. Abdulmalik ,
Opeyemi A. Abisoye , and Sulaimon A. Bashir 

Department of Computer Science, Federal University of Technology, Minna, Nigeria
bennychika2@gmail.com, {drmalik,o.a.abisoye,
bashirsulaimon}@futminna.edu.ng

Abstract. Data mining technique has been used to extract potentially useful knowledge from big data. However, data mining sometimes faces the issue of incorrect results which could be due to the presence of an outlier in the analyzed data. In the literature, it has been identified that the detection of this outlier could enhance the quality of the dataset. An important type of data that requires outlier detection for accurate prediction and enhanced decision making is time series data. Time series data are valuable as it helps to understand the past behavior which is helpful for future predictions hence, it is important to detect the presence of outliers in time series dataset. This paper proposes an algorithm for outlier detection in Multivariate Time Series (MTS) data based on a fusion of K-medoid, Standard Euclidean Distance (SED), and Z-score. Apart from SED, experiments were also performed on two other distance metrics which are City Block and Euclidean Distance. Z-score performance was compared to that of inter-quartile. However, the result obtained showed that the Z-score technique produced a better outlier detection result of 0.9978 F-measure as compared to inter-quartile of 0.8571 F-measure. Furthermore, SED performed better when combined with both Z-score and inter-quartile than City Block and Euclidean Distance.

Keywords: Outlier detection · Time series data · Multivariate · Outliers · K-Medoid · Euclidean distance · Z-scores · City block

1 Introduction

Currently, the volume of data is significantly increasing with each passing day in various application fields, hence, it is essential to use more effective technology to examine and manage these volumes of data to discover the integral, previously unknown, and potentially useful knowledge. One key data processing technology that is used to examine and manage data is data mining. Data mining is a method of obtaining valuable information from unprocessed-data [1]. With data mining, a search engine can be used to analyze massive quantities of data and detect interesting trends automatically without needing human involvement [2]. Data mining is generally about finding non-trivial, concealed,

and important information from different data types. There are many fields in which data mining is commonly used: retail, network security, financial data processing, biochemical data analysis, telecommunications industry, and other research applications [1].

One of the fundamental issues of data mining is to achieve a profound and accurate conclusion of the study. And outlier detection is a crucial step in solving this basic problem [3, 4]. An outlier is a data point that does not fit into the usual data set classification points. Outlier detection has major applications in data preprocessing and the elimination of abnormal points for stock price prediction, credit scoring, advertisement, intrusion detection, system evaluation, and e-mail spam identification among others [5]. Outlier detection is also known in some researches as anomaly detection and is a critical longstanding research concern in the fields of data mining and statistics. The principal principle behind finding an outlier is to locate unusual points in the data points. Outlier detection points out artifacts that often deviate from a given set of data [6] and Outliers may potentially contain some valuable information [7]. For example, the identification of unusual trends in patient databases can be helpful for disease identification in clinical applications [5, 8].

Time series data is an essential form of data that needs outlier detection for accurate prediction and effective decision making. Data from time series are valuable because it helps to explain the past actions that would be useful for future predictions [9, 10]. Time series data can be used for business forecasting which in turn can improve business decision making. These data are applicable in various fields such as mathematics, sensor processing, pattern recognition, economics, mathematical finance, weather monitoring, earthquake prediction, electromyography, industrial engineering, astronomy, telecommunications and, in large part, any field of applied engineering and science [11]. The identification of outliers in time series data is therefore critical for better decision making and prediction. Considering the importance of outlier detection in time series data, numerous researches have been carried out using various types of algorithms in the success of an outlier detection function in time series data. Algorithms that have been proposed in the literature for outlier detection in multivariate time series dataset include genetic algorithm, Mahalanobis distance, Euclidean distance, autoregressive integrated moving average with exogenous inputs (ARIMAX) model and k-means algorithm while z-score and inter-quartile (box plot) algorithm have been used for univariate time series dataset. This paper proposes a method that combines K-Medoid which is a cluster-based algorithm, Standard Euclidean Distance which is a distance-based, and Z-Score algorithm for outlier detection in multivariate time series (MTS) data. This study's principal contributions are as follows:

1. A fusion of cluster-based and distance-based algorithms for improved outlier detection.
2. Comparative experimentation of different distance-based algorithm on the data cluster obtained from the clustering algorithm.
3. Extension of Z-score and Box plot algorithms for effective outlier detection in MTS dataset.

The organization of these studies is as follows: section two presents the relevant work, section three presents the methods used, section four describes the findings and discussion and finally, in section five and section six conclusions were drawn and recommendations for future works were presented.

2 Related Work

The presence of outliers in time series and non-time series data could provide important information that can be used by an analyst in drawing conclusions. Also, an outlier presence could alter the result produced during analysis. Hence several types of research have been carried to detect these outliers and even remove them when necessary.

Jones et al. [12] performed outlier detection in physical activity using the k-means clustering algorithm. In this work, a FilterK technique was introduced to boost the quality of clusters of k-means based on physical activity. The FilterK algorithm uses an outlier score function that allocates a level of abnormality based mostly on the standardized score obtained from each of its 3 tests which are average distance to neighbors, distance to the closest centroid, and density of the neighborhood. This allows the setting of a standard outlier threshold score as an exclusion criterion when searching the data for abnormalities on the accelerometer. However, the robustness, strength, and weakness of this proposed method were not tested as the method was tested on a small range of accelerometer dataset.

Souza, Aquino and Gomes [13] integrated tensor decomposition with data categorization to outliers in urban spaces detect applications and provide useful information for urban planning and operation. This method consists of three stages which are dimensionality reduction stage in which several latent variables are derived from the contraction, classification of latent variables stage in which the latent variables produced are being used for classification to obtain high-quality groups from the factorization stage, and finally, the production of a polished environmental sequence identification phase that deals with the design of a monitoring process statistics to detect events outside the regularity trends of the measured dataset. A drawback of this proposed approach is that it needs a large sample size to classify the outliers and this necessity has an effect on applications in real-time.

Erkus and Purutic [4] proposed a non-parametric method for quasi-periodic outliers detection in time series data based on a combination of frequency-domain and Fourier transform definition. The result of the algorithm was compared to four existing algorithms namely: Grubbs, box-plot, autonomous anomaly detection, and z-score method, and it was proven that the developed approach worked even better if the data had quasi-periodic structures from start to finish. However, this method did not consider outlier detection for high dimensional datasets as only outliers in univariate dimensions were considered.

Wang et al. [14] considered outliers detection based on its importance in air pollution forecasting. It was found that the identification and correction of the outlier point of the original time series affect air pollution prediction. To improve their air quality index forecasting, a novel hybrid approach based on outlier detection, corrective algorithm, and a heuristic intelligent optimizing algorithm was developed. To achieve outlier identification and correction, the Hampel recognition system was based on measured local

median and standard deviation was adopted. Only the air quality index time series was considered while other influencing factors were not considered in this work.

Borah and Nath [15] proposed a method for solving the problem of outlier detection for incremental medical data, as most of the current outlier recognition algorithms are capable of controlling just static data, and must therefore perform incremental data from scratch. Rare pattern-based outlier detection (RPOD) technique was used to carry out this task. The RPOD consists of two stages: the full info on the database is stored in a compact prefix-tree structure at the first stage. So then the desired set of uncommon patterns is extracted progressively from the prefix-tree by skimming down the index only once. While the second phase includes the identification of outliers based on the unusual patterns obtained during the first stage. The outliers were identified using 3 techniques of outlier detection: Rare Pattern Support Deviation Factor (RPSDF), Transaction Outlier Factor (TOF), and Rare Pattern Outlier Factor (RPOF), respectively. The drawback of this approach is that by setting a minimum support level, RPOD creates outliers based upon the notion of uncommon patterns. And the number of uncommon patterns produced is basically too large and it can be very cumbersome to assign an acceptable support value for outlier detection

Ghallab, Fahmy and Nasr [16] model called NRDD-DBSCAN focused on density-based spatial clustering of noise-based applications (DBSCAN) algorithms and using resilient distributed data sets (RDDs) to identify anomalies that affect Internet of Things (IoT) technology data quality. The NRDD-DBSCAN has been used to solve the RDD-DBSCAN system's low-dimensionality problem and also to solve the DBSCAN algorithm issue not being able to handle IoT data. And NRDD-DBSCAN can also be used to boost the efficiency of the present data in IoT devices and applications. The primary drawback of this proposed system is that the data reduction and distribution process uses principal component analysis (PCA) to reduce the dimensionality of the datasets, however, PCA works well by reducing the dimensionality of a strongly correlated linear n-dimension data that prevents the system from processing non-linear data.

3 Methodology

This section discusses the techniques that were used to carry out this research work. Figure 1 shows the processes and techniques that were used to achieve the research goal. Each of the steps shown in the diagram in Fig. 1 is discussed in the sub-sections below.

3.1 Data Acquisition

The time series dataset used in this study was gotten from the Time series section in the UCI repository. Two labeled multivariate time series datasets were used which are the occupancy dataset and the MHEALTH dataset. The MHEALTH which is a multivariate time series dataset consists of observations of body movement and vital signs for 10 individuals of different personalities while undertaking 12 physical activities, and it consists of 23 feature vectors. The occupancy data set defines the dimensions of a room and is intended to determine whether space is occupied or not and it has 20,560 one minute observations taken over a period of few weeks.

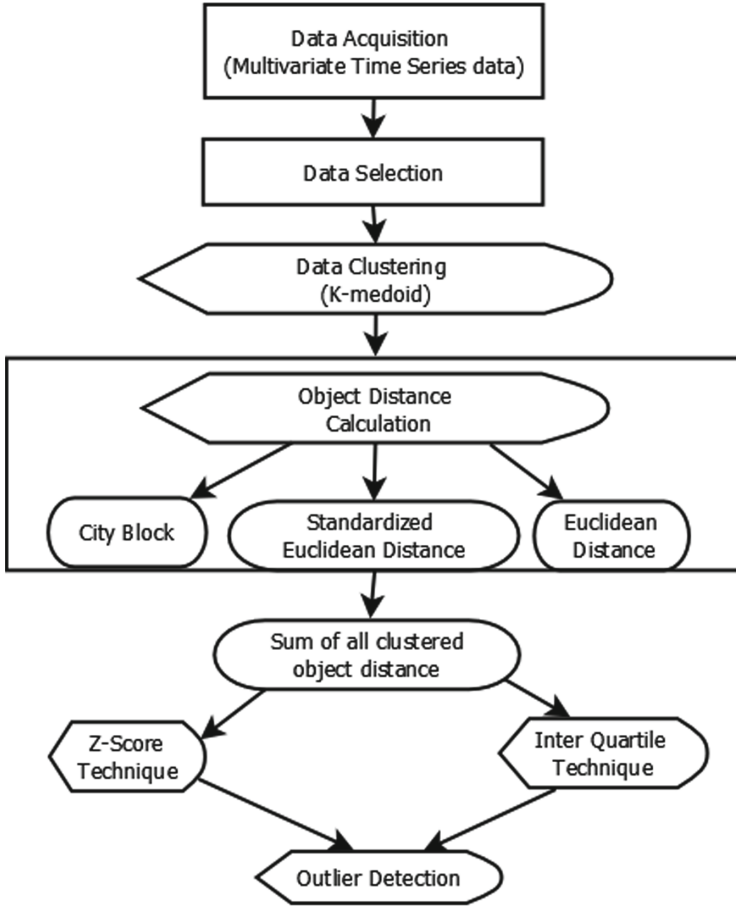


Fig. 1. Proposed system

3.2 Data Selection

In this work, not all instances in the dataset were used and this was done in order to make the dataset suitable for the outlier detection task. For the MHealth dataset 2000 instances were used out of 161280 instances contained in the MHealth dataset. Out of the 2000 instances used, 1975 instances are labeled as a normal class while 25 instances were inserted as outliers. For the occupancy dataset, 1500 observations were taken out of 9753 observations contained in the dataset and the 1500 observations used consist of 1485 occupied rooms observations labeled as a normal class and 15 empty room observations identified as outliers.

3.3 Data Clustering (K-Medoid)

In the literature, the clustering result of the k-medoid clustering method cannot be influenced by the presence of outliers as data points are selected to be the medoids. K-medoid

algorithm is a partitioning clustering algorithm that divides a set of n objects into k number of clusters. This approach is widely used in areas that demand robustness to outlier data, arbitrary distance measures, or those that do not have a precise definition of mean and median [17]. Given objects with p variables that will be grouped into k ($k < n$) clusters. Defining the j th vector of object i as Y_{ij} ($i = 1, \dots, n; j = 1, \dots, p$). K-medoid algorithm consists of three steps which are:

1. Selecting initial medoids: using Euclidean distance as a metric of dissimilarity, compute the difference between each pair of objects using the formula below:

$$D_{ij} = \sqrt{\sum_{a=1}^p (Y_{ia} - Y_{ja})^2} \quad i = 1, \dots, n; j = 1, \dots, n \quad (1)$$

2. Calculate p_{ij} by using the formula in Eq. 2 below to make an initial guess at the cluster centers.

$$P_{ij} = \frac{d_{ij}}{\sum_{i=1}^n d_{ij}} \quad i = 1, \dots, n; j = 1, \dots, n \quad (2)$$

3. Replace current medoid with the object in each cluster which minimizes the distance to other objects in its cluster.

In this study, k was assigned a value of three (3).

3.4 Object Distance

After grouping data into clusters using the k-medoid algorithm, the distance of each object from the medoid of each cluster was calculated using three distance-based algorithms which are: the SED, City Block, and the Euclidean distance. Considering the m -by- n data matrix X , which will be handled as m (1-by- n) row vectors x_1, x_2, \dots, x_m , the different distances between both the x_s and x_t vectors are stated as follows:

1. **Euclidean Distance:** A linear distance in Euclidean space between two points [18].

$$D_{st}^2 = (x_s - x_t) (x_s - x_t)' \quad (3)$$

2. **Standardized Euclidean Distance (SED):**

$$D_{st}^2 = (x_s - x_t) U^{-1} (x_s - x_t)' \quad (4)$$

Where U is the n -by- n diagonal matrix with diagonal elements given by s_j^2 , which signifies the variance of the variable x_j over the m objects.

3. **City Block Distance:** the city block distance between two vectors x_s and x_t in an n -dimensional vector space with static Cartesian coordinate is the addition of the lengths of the projections of the line section between the points onto the coordinate axes. The formula is as follows:

$$D_{st} = \sum_{j=1}^n |x_{sj} - x_{tj}| \tag{5}$$

3.5 Distance Summation

After using the above distance metrics to calculate the distance of objects from each cluster medoid, this calculated distances for all clusters are summed together for each object. Hence a single distance gotten from the summation of all an object distance from each cluster medoid is used as an input for the inter-quartile and z-score technique. Given that there are i ($i = 1, \dots, n$) instances grouped into 3 clusters with various calculated distance (C_{i1}, C_{i2}, C_{i3}). This distance summation is given as:

$$d_s = C_{i1} + C_{i2} + C_{i3} \tag{6}$$

3.6 Outlier Detection Techniques

Two techniques were considered for outlier detection and these techniques are the z-core and inter-quartile techniques.

1. **Z-core technique:** This is a statistical measurement of a score’s correlation to the mean in a collection of scores [19]. A Z-score of zero indicates that the score is the same as the average. It can also be a negative/positive value, signifying if it is below or above the average and by how many standard deviations [20]. Z-score makes use of average and standard deviation values to detect outliers in a dataset. Z-score is calculated as follows:

$$Z_{score}(i) = \frac{x_i - \mu}{sd} \tag{7}$$

Where μ = distribution mean, sd = standard deviation and x_i = each object in the distribution. The standard deviation can be calculated using:

$$sd = \sqrt{\frac{1}{n - 1} \sum_{i=1}^n (x_i - \mu)^2} \tag{8}$$

To detect outliers based on z-score a cut-off value of 3 was used. That is any distance value with an integer value greater than 3 are detected as outliers.

2. **Inter-quartile technique:** this technique is also known as box plot method. It is called the box plot technique because it forms a box chart, which shows information such as the upper and lower outliers of the univariate data set, lower and upper quartiles, and the median. This method has only been used to detect outliers in a univariate dataset but this study extends its usage to a multivariate dataset. This extension is possible as the sum of the distances was used as a single univariate input to the box-plot and based on input the technique was able to detect outliers in a multivariate dataset. The rule for this method is as follows.

- a. Calculate the interquartile range (IQR)

$$\text{IQR} = Q_3 - Q_1 \quad (9)$$

- b. Calculate the lower and upper internal boundaries using the formulas:

$$\text{LIB} = Q_1 - 1.5\text{IQR} \quad (10)$$

$$\text{UIB} = Q_3 - 1.5\text{IQR} \quad (11)$$

- c. Calculate the lower and upper external boundaries using the formulas.

$$\text{LEB} = Q_1 - 3\text{IQR} \quad \text{Q}_{-1} - 3\text{IQR} \quad (12)$$

$$\text{UEB} = Q_3 - 3\text{IQR} \quad (13)$$

- d. Observation values between internal and external boundaries are defined as possible outliers.
 e. Observation values outside the external boundaries are determined as outliers.

3.7 Performance Metrics

Four performance measures were used to evaluate the proposed method. These measures are explained below.

- **Precision:** is a measure that evaluates the number of accurate predictions made correctly. It is determined as the ratio of correctly predicted positive examples, divided by the total number of predicted positive examples.

$$\text{Precision} = \text{True_Positives} / (\text{True_Positives} + \text{False_Positives}) \quad (14)$$

- **Recall:** is a statistic that evaluates the amount of accurate positive predictions that could have been made from all the positive predictions.

$$\text{Recall} = \frac{\text{True_Positives}}{\text{True_Positives} + \text{False_Negatives}} \quad (\text{SEQ "equation" } \backslash n \backslash * \text{ MERGEFORMAT 15}) \quad (15)$$

- **F-measure:** It is the harmonious measure of precision and the recall

$$\text{F - measure} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \quad (\text{SEQ "equation" } \backslash n \backslash * \text{ MERGEFORMAT 16}) \quad (16)$$

- **Accuracy:** Accuracy is basically defined as the rate of correct classifications. And it is calculated as follows:

$$\text{Accuracy} = \frac{\text{True Positive} + \text{True negative}}{\text{True Positive} + \text{True negative} + \text{False Positive} + \text{False negative}} \quad (\text{SEQ "equation" } \backslash n \backslash * \text{ MERGEFORMAT 17}) \quad (17)$$

This study deals with imbalanced classification problems where the number of examples in the dataset for each class label is not balanced, that is the distribution of examples across the known class is skewed. For example, in the occupancy dataset, only 15 examples are of class 1 (outliers) while the remaining 904 examples are class 0 (normal). Hence the most important performance metrics considered for evaluation are precision, recall, and f-measure which evaluates imbalanced classification problems effectively.

3.8 Algorithm Fusion

After performing the data acquisition and data selection process. The K-medoid, Standard Euclidean Distance (SED), and Z-score were combined as follows: The data were then grouped into 3 clusters using the k-medoid algorithm. After clustering the distance of each object from the medoid of each cluster was calculated using Standard Euclidean Distance. Given that there are 3 clusters 3 different SED distances were obtained for each object from each cluster. In order to utilize the 3 SED distances obtained for each object a single distance for each object was generated by summing all of the 3 cluster distances. This generated single distance was used as input data to the z-score technique. The z-score of the SED distances was computed and based on a threshold value of the z-score a data was identified as an outlier.

4 Results and Discussion

In this study, outlier detection was performed using K-Medoid, SED, and Z-Score. The performance of SED was compared to that of Euclidean Distance and City clock. And the performance of z-score was compared to that of inter-quartile. The algorithms were tested using the MHealth and Occupancy multivariate time series dataset. The results obtained are presented in Tables 1, 2, 3, and 4 below. Table 1 and Table 2 shows the results obtained for MHealth multivariate time series dataset.

Table 1. Results of inter-quartile outlier detection technique for the three different distance metric for MHealth dataset

Outlier detection algorithm	Precision	Recall	F-measure	Accuracy
K-Medoid + SED + inter-quartile	0.2778	1.0000	0.4348	0.9133
K-Medoid + Euclidean + inter-quartile	0.2400	0.933	0.3889	0.9022
K-Medoid + City Block + inter-quartile	0.2727	1.0000	0.4286	0.9111

The results obtained for inter-quartile outlier detection technique based on the three different distance metric is shown in Table 1 above. Table 1 shows that standardized Euclidean distance performed better with an f-measure of 0.4348 as compared to Euclidean and City block distance with f-measure of 0.3889 and 0.4286 respectively.

Table 2. Results of Z-score outlier detection technique for the three different distance metric for MHealth dataset

Outlier detection algorithm	Precision	Recall	F-measure	Accuracy
K-Medoid + SED + Z-score (Proposed System)	0.9704	0.9793	0.9748	0.9511
K-Medoid + Euclidean + Z-score	0.9748	0.9770	0.9759	0.9533
K-Medoid + City Block + Z-score	0.9725	0.9747	0.9736	0.9489

Table 2 shows the results obtained for the z-score detection technique based on the three different distance metrics for the MHealth dataset. Based on the precision, recall, and f-measure obtained in Table 2 it can be seen that the z-score technique performs better with f-measures of 0.9748, 0.9759, and 0.9736 for Standardized Euclidean, Euclidean, and city block respectively than the inter-quartile technique with f-measure of 0.4348, 0.3889 and 0.4286 for Standardized Euclidean, Euclidean and city block respectively. Table 3 and Table 4 shows the results obtained for Occupancy multivariate time series dataset.

Table 3. Results of inter-quartile outlier detection technique for the three different distance metric for Occupancy dataset

Outlier detection algorithm	Precision	Recall	F-Measure	Accuracy
K-Medoid + SED + inter-quartile	0.7500	1.0000	0.8571	0.9946
K-Medoid + Euclidean + inter-quartile	0.4286	1.0000	0.6000	0.9782
K-Medoid + City Block + inter-quartile	0.7500	1.000	0.8571	0.9946

Table 3 shows the results obtained for the inter-quartile outlier detection technique based on the three different distance metrics for the occupancy dataset. From Table 3 it can be seen that standardized Euclidean distance and city block distance performed better with an f-measure of 0. 8571 as compared to Euclidean with an f-measure of 0. 6000.

Table 4. Results of Z-score outlier detection technique for the three different distance metric for Occupancy dataset

Outlier detection algorithm	Precision	Recall	F-Measure	Accuracy
K-Medoid + SED + Z-score (Proposed System)	1.000	0.9956	0.9978	0.9956
K-Medoid + Euclidean + Z-score	1.000	0.9934	0.9967	0.9935
K-Medoid + City Block + Z-score	1.000	0.9956	0.9978	0.9956

Table 4 shows the results obtained for z-score detection technique based on the three different distance metrics for the Occupancy dataset. Based on the precision, recall, and f-measure obtained in Table 4 it can be seen that the z-score technique still performs better with f-measures of 0.9978, 0.9967, and 0.9978 for Standardized Euclidean, Euclidean and city block respectively than the inter-quartile technique with f-measure of 0.8571, 0.6000, and 0.8571 for Standardized Euclidean, Euclidean, and city block respectively.

From the results obtained for the inter-quartile technique for both occupancy and MHealth dataset as shown in Table 3 and Table 1 respectively, it can be seen that squared Euclidean distance performs better than Euclidean and city block distance. Hence it can be concluded that the inter-quartile technique performs better when combined with squared Euclidean distance. Also from Table 2 and Table 4, it can be seen from the precision, recall, and f-measure that all three distance metrics produce good and satisfactory results when combined with z-score. However, in both Mhealth and occupancy dataset z-score produces a better result as compared to inter-quartile technique. With high precision, recall, f-measure, and accuracy it can be deduced that outliers can be detected effectively with z-score when combined with any of the distance metric especially with the Standardized Euclidean distance. Hence a system which combines k-medoid, SED and z-score is proposed for effective outlier detection in multivariate time series dataset.

5 Conclusion

This study was able to extend the inter-quartile and z-score technique for outlier detection on multivariate time series dataset as compared to existing works which used these techniques for outlier detection on only univariate datasets. This is due to the ability of the algorithm to combine the multivariate features into a univariate feature. Also, this study was able to detect outliers more effectively as compared to existing works on outlier detection. From the study, it can be concluded that the z-score technique performs better than the inter-quartile technique for outlier detection.

6 Future Works

In this study, only the extreme outlier values generated by the inter-quartile techniques were identified as outliers while the mild identified outlier values were ignored and assigned to the normal class because the data used are skewed it caused a large number of observations to be determined as outliers in the data. This is due to the use of the interquartile distances and lower and upper quadrants measured without considering the skewness of the data set. Hence is it recommended that the skewness of the dataset should be considered to improve the performance of the inter-quartile technique. Three distance metric was used in this work, for future work other distance metrics can be used. This work can be extended for outlier detection for other non-time series datasets.

References

1. Jain, S., Sahib, F., Kaur, A., Sahib, F.: A review paper on comparison of clustering algorithms based on outliers, vol. 3, no. 05, pp. 178–182 (2016)
2. Rajagopal, S.: Customer data clustering using data mining technique, vol. 3, no. 4 (2011). <https://doi.org/10.5121/ijdms.2011.3401>
3. Ramesh, K.B., Aljinu, K.K.V.: A survey on outlier detection techniques in dynamic data stream. *Int. J. Latest Eng. Manag. Res. IJLEMR* **02**(08), 23–30 (2017)
4. Erkus, E.C., Purutç, V.: Journal Pre-proof. *Eur. J. Oper. Res.* (2020). <https://doi.org/10.1016/j.ejor.2020.01.014>
5. Pamula, R., Deka, J.K., Nandi, S.: An outlier detection method based on clustering. In: *Proceedings - 2nd International Conference on Emerging Applications of Information Technology, EAIT 2011*, pp. 253–256, February 2011. <https://doi.org/10.1109/eait.2011.25>
6. Akouemo, H.N., Povinelli, R.J.: Time series outlier detection and imputation, pp. 1–5 (2014)
7. Jiadong, R., Hongna, L., Changzhen, H., Haitao, H.: ODMC: outlier detection on multivariate time series data based on clustering. *J. Converg. Inf. Technol.* **6**(2), 70–77 (2011). <https://doi.org/10.4156/jcit.vol6.issue2.8>
8. Ren, J., Li, H., Hu, C., He, H.: ODMC: outlier detection on multivariate time series data based on clustering. *J. Converg. Inf. Technol.* **6**(2), 70–77 (2011). <https://doi.org/10.4156/jcit.vol6>
9. Liu, S., Wright, A., Hauskrecht, M.: Online conditional outlier detection in nonstationary time series. In: *Association Advance Artificial Intelligence* (2017)
10. Abayomi-Alli, A., Odusami, M.O., Abayomi-Alli, O., Misra, S., Ibeh, G.F.: Long short-term memory model for time series prediction and forecast of solar radiation and other weather parameters. In: *2019 19th International Conference on Computational Science and Its Applications (ICCSA)*, Saint Petersburg, Russia, pp. 82–92, July 2019. <https://doi.org/10.1109/iccsa.2019.00004>

11. Hasan, E.A.: A method for detection of outliers in time series data. *Int. J. Chem. Math. Phys. IJCMP* **3**(3), 56–66 (2019). <https://doi.org/10.22161/ijcmp.3.3.2>
12. Jones, P.J., et al.: FilterK : a new outlier detection method for k-means clustering of physical activity. *J. Biomed. Inform.* **104**, 103397 (2020). <https://doi.org/10.1016/j.jbi.2020.103397>
13. Souza, T.I.A., Aquino, A.L.L., Gomes, D.G.: A method to detect data outliers from smart urban spaces via multiway analysis. *Future Gener. Comput. Syst.* (2018). <https://doi.org/10.1016/j.future.2018.09.062>
14. Wang, J., Du, P., Hao, Y., Ma, X., Niu, T., Yang, W.: An innovative hybrid model based on outlier detection and correction algorithm and heuristic intelligent optimization algorithm for daily air quality index forecasting. *J. Environ. Manage.* **255**, 109855 (2020). <https://doi.org/10.1016/j.jenvman.2019.109855>
15. Borah, A., Nath, B.: Journal pre-proof. *Appl. Soft Comput. J.* 1–51 (2019). <https://doi.org/10.1016/j.asoc.2019.105824>
16. Ghallab, H., Fahmy, H., Nasr, M.: Detection outliers on Internet of Things using big data technology. *Egypt. Inform. J.* **21**, 1–8 (2019). <https://doi.org/10.1016/j.eij.2019.12.001>
17. Hudaib, A., Khanafseh, M., Surakhi, O.: An improved version of K-medoid algorithm using CRO. *Mod. Appl. Sci.* **12**(2), 116–127 (2018). <https://doi.org/10.5539/mas.v12n2p116>
18. Dokmani, I., Parhizkar, R., Ranieri, J., Vetterli, M.: Essential theory, algorithms and applications. *IEEE Signal Process. Mag.* **32**, 1–17 (2015)
19. Kolbaşı, A., Ünsal, P.A.: A comparison of the outlier detecting methods: an application on Turkish foreign trade data. *J. Math. Stat. Sci.* **5**, 213–234 (2015)
20. Anuradha, C., Murty, P.S.R.C., Kiran, C.S.: Detecting outliers in high dimensional data sets using Z-score methodology. *Int. J. Innov. Technol. Explor. Eng. IJITEE* **9**(1), 48–53 (2019). <https://doi.org/10.35940/ijitee.a3910.119119>