

## SMART OUTLIER DETECTION OF WIRELESS SENSOR NETWORK

Sahar Kamal<sup>1</sup>, Rabie A. Ramadan<sup>2</sup>, Fawzy EL-Refai<sup>3</sup>

<sup>1</sup>Department of Electronics and Electrical Communications,  
Higher Institute of Engineering, El-Shorouk Academy, El-Shorouk city, Egypt

<sup>2</sup>Computer Engineering Department, Cairo University, Egypt

<sup>3</sup>Department of System and Computer Engineering, El-AZhar University, Cairo, Egypt

**Abstract.** *Data sets collected from wireless sensor networks (WSN) are usually considered unreliable and subject to errors due to limited sensor capabilities and hard environment resulting in a subset of the sensors data called outlier data. This paper proposes a technique to detect outlier data base on spatial-temporal similarity among data collected by geographically distributed sensors. The proposed technique is able to identify an abnormal subset of data collected by sensor node as outlier data. Moreover, the proposed technique is able to classify this abnormal observation, an error data set or event affected set. Simulation result shows that high detection rate is achieved compared to conventional outlier detection techniques while preserving low positive false alarm rate.*

**Key words:** *wireless sensor network, outlier's detection, fuzzy logic, spatial and temporal similarity*

### 1. INTRODUCTION

Wireless sensor network is considered a promising solution for monitoring and measurement of natural physical phenomena such as temperature, humidity, earthquakes, pressure, light, volt, etc. A typical WSN consists of a large number of very small sensors deployed over a topological area of interest. These sensors are supplied by power resources (batteries, solar cells), measurement unites, processing units and wirelesses TX/RX unit. Unfortunately, the data collected from sensor nodes are considered inaccurate and may be even unreliable due to measurement errors or superimposed noise on the received data packets in [2]. Duplicated measurement or even missing values are not common in the data set collected by a WSN. A subset of data which appear to be in consistence with the whole

---

Received August 30, 2015; received in revised form November 15, 2015

**Corresponding author:** Rabie A. Ramadan

Computer Engineering Department, Cairo University, Egypt

(e-mail: rabie@rabieramadan.org)

\*An earlier version of this paper was presented at the International Conference on Recent Advances in Computer Systems RACS-2015, Hail University, Saudi Arabia, 2015 [1].

data set from which it is collected is called an outlier. Outlier can be defined as in [2] "an outlier is a subset of observations which appear to be inconsistent with other dataset". On the other hand, outliers as in [3] can be defined as "those measurements that are deviated from the consistence dataset". Each of two definitions can be used as a solution to declare the outlier in a data set. Abrupt events such as sudden sensor failure, battery power deployment or even natural physical phenomena are also reasons to which outlier data can be attributed. In order to boost the accuracy and reliability of the collected sensor data, an outlier detection process should be applied and possibly corrected.

There are three sources of outliers due to environmental changes or error coming from a faulty sensor, which can be defined as (1) errors& noise, (2) events and (3) malicious attacks, the last one being related to the network security as in [2]. Noise or error refers to a noise-related measurement or data instance coming from a faulty sensor. Outliers caused by errors may occur frequently, while outliers caused by events tend to have a smaller probability of occurrence. Erroneous data is normally represented as an arbitrary change and is extremely different from the rest of the data. Noisy data as well as erroneous data should be eliminated or corrected if possible. However, events may arise due to sudden change in the real world, for example rainfall, forest fire, chemical spill, air pollution, etc. Removing the event outlier from data set will lead to a loss of important hidden information of the data about events as in [4]. Outliers that are very close to random errors in terms of size can only be determined through the application of outlier tests. Outlier classification as an event or error is an important matter. Many researches consider outliers and events as similar conditions by treating events as some sort of outliers. Due to the fact that there are spatial-temporal similarities between neighboring nodes, measurements enable us to classify outlier as either an event or error. This depends on the fact that error data observations seem to be unrelated, while event observations seem to be spatially correlated as in [5].

The main approaches to determine outliers can be grouped as Statistics-based methods, Nearest Neighbor-Based, Cluster-Based and Artificial Intelligence techniques. New approaches are used for outlier detection including Artificial Intelligence techniques such as Neural Networks and Fuzzy Logic technique. The latter was suggested by [6] in which it can also be used for geodetic networks for outlier detection. The main aim of outlier detection in WSN is to declare outliers with high detection rate while decreasing the resource consumption of network.

Our work is based on the observation that in most applications of WSNs measurements of sensors in the environment tend to be highly correlated for sensors that are geographically close to each other (spatial similarity), and also highly correlated for a period of time (temporal similarity) as in [5]. Using this observation, we take advantage of the spatial and temporal similarity in the sensor data. In the first study, we detect outliers in the univariate attribute in WSN. The main contribution of this paper is the use of Euclidean distance and fuzzy logic to detect outliers in wireless sensor networks. However, spatial and temporal similarity were used to make it easy to distinguish between error and event. If probability of output of fuzzy logic is above a prefixed threshold, the observation is considered as an outlier. The model is tested on a real data set from Grand-St-Bernard as in [7] and implemented using MATLAB. This paper achieves a high detection rate and still keeps a low false positive alarm rate and computational complexity.

The rest of the paper is organized as follows: Section (2) shows the necessary background definition related to outlier detection. The proposed algorithm is presented in section (3) along with the assumptions upon which the proposed technique is built. Section (4) shows

experimental results and the performance evaluation of the proposed technique using a realistic data set. Finally, the whole paper is concluded in section (5).

## 2. RELATED WORK

Recently, there are many researches in outlier detection of WSN to improve reliability and quality of measurement sensor. These researches used different techniques to detect outlier such as statistical-based, nearest neighbor-based, clustering-based, classification-based, and spectral decomposition-based approaches. In general, these researches can be those that do not use spatial or temporal correlation data set or those that are based on spatial or temporal correlation only or on both. In 2006, the author in [8], uses the spatial correlation that exists among neighboring sensor nodes to distinguish between outlying sensors and event boundary. In this model, each node calculates the difference between its own measurements and the median from its neighboring measurements. Then outlying node is declared when the absolute value of its measurement's deviation degree is greater than a pre-selected threshold. This technique suffers from a low detection rate because it ignores temporal correlation between sensor data reading. As shown by [9], this model used a cluster based technique to identify the global outlier. First, each node clusters the reading and reports cluster summaries and then transmits the raw sensor reading to its cluster head. The cluster head collects cluster summaries from all of its nodes before sending them to the sink. An outlier cluster can be declared in the sink if the cluster's average inter-cluster distance is greater than one threshold value of the set of inter-cluster distances. However, these models suffer from the choice cluster width parameter. Additionally, these techniques increase computational complexity when computing the distance between data instance. In [10] author uses distance similarly to identify global outliers in WSN. Each node uses a distance in a similar way to identify local outliers and then broadcast abnormal data Instances to all neighboring node for verification. This technique is repeated until all neighboring nodes agree on the global outliers. This technique increases computational complexity and it isn't adapted for a large scale network. In 2007, the proposed technique as in [11] uses one class quarter sphere based technique to detect outliers in WSN. This technique takes advantage of temporal correlation to identify local outliers at each node. A measurements sensor that lies outside the quarter sphere is considered as an outlier. Each node transmits only brief information to its parent for global outlier's classification. This technique suffers from a low detection rate because it ignored spatial correlation between neighboring nodes. At 2008, the author as in [12] uses a centered quarter-sphere support vector to detect local outlier in WSN. This technique takes advantage of spatial correlations that exist in sensor data of adjacent nodes to reduce the false alarm rate and to distinguish between events and errors, but it ignores temporal correlation and increases computational complexity. But in 2009, the author as in [13] used outlier detection technique to identify outliers in data set of WSN. This technique takes advantage of spatial temporal correlation exist among sensor data reading. In 2011, author as in [14] proposed outlier detection method in the wireless sensor networks and distinguishes between event and error. This technique is used to classify the sensor node data as local outlier or cluster outlier or network outlier. This technique considers the network outlier or cluster outlier as event and local outlier as error. This algorithm suffers from high computational complexity. In 2012, the author of [15] use the advantage of temporal correlation only to detect the outlier in WSN. However, this technique suffers from some computational complexity. This approach

differs from our approach in that our approach has the advantage of spatial-temporal similarity combined with fuzzy logic to detect outlier and identify errors and events with high detection rate and relatively low false positive rate in comparison with the result in [15]. In 2013, the author as in [16] uses temporal and spatial properties to identify outliers and distinguish between event and error but with low detection rate and false positive rate in comparison with our approach.

### 3. THE PROPOSED STODM TECHNIQUE

Sensor nodes are assumed to be densely deployed and synchronized in WSN. A subset of sensors is considered as members of the same cluster if they fall within the same radio transmission range of each other. At any time interval  $\Delta t$ , each node reads a data vector  $s_{ij}$  where “ $i$ ” is the time index of the data symbol and “ $j$ ” is the node spatial ID. The potential of an outlier detection technique is to identify a subset  $x_i$  of each sensor set  $s_i$  as outliers. A super advantage of a given detection technique is to classify deviation data instance as event or error.

In this section, the proposed approach is introduced in details. Many outlier detection techniques have been developed, however, they did not take into account the interesting events. On the other hand, several recently developed researches are interested only in events and did not care about erroneous data. In this paper, a new distance-based approach depends on spatial-temporal similarity combined with fuzzy logic-based approach is proposed to classify outliers, i.e. error data or events. Our methodology consists of the following steps: first step the spatial and temporal similarity is calculated, each one of these is entered as input or (membership function) to fuzzy logic to detect outliers in each node. Second step classifies the outlier as event or error.

#### 3.1. Spatial-temporal similarity

In our proposed algorithm, spatial-temporal similarity is calculated using a two-step process.

First step, the temporal similarity of a given data set of sensor node is calculated on point by point basis and is given by first order difference  $|s_{i2}-s_{i1}|$ . The absolute difference is compared to a pre-specified threshold which is calculated according to tolerance of temperature sensor. A data point  $s_{i2}$  is considered similar to other points if the absolute first order difference does not exceed the threshold. Otherwise, dissimilarity is obtained and point of data may be outlier.

Second step, spatial similarity is calculated based on the distance between neighboring nodes. We use the Euclidean distance to calculate similarity measure between two points -  $x$ ,  $y$ , that are in the same transmission range and are in the same close time which is calculated as Eq. (1). Euclidean Distance is a popular choice for univariate and multivariate continuous attributes as [17]. Data instance in point  $x$  is considered similar to data point in  $y$  if Euclidean distance  $d(x, y)$  does not exceed Preselected threshold. Spatial link is defined as number of spatial similarity to each point with its neighbors as in Eq. (2). Where spatial similarity threshold is calculated by computing mean distance of all data points in the close time.

$$D(x, y) = \sqrt{(x - y)^2} \quad (1)$$

$$\text{Spatial link} = \sum_{i=1}^n no_{of} \text{similarity to each node} \quad (2)$$

where  $n$  is the number of neighboring nodes.

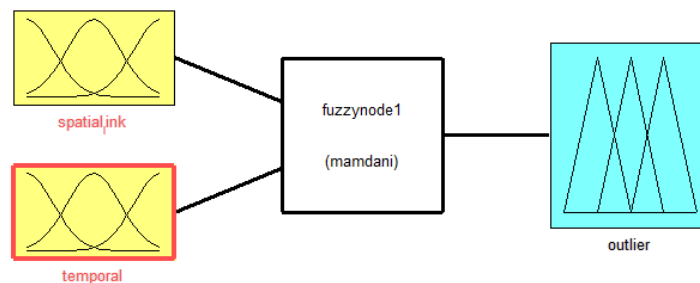
### 3.2. Fuzzy logic model

Recently, many approaches have been tested on decision making theories. Some of the artificial techniques that are used in outlier analysis are Neural Networks, Support Vector Machine and Fuzzy Logic as in [18]. Our approach use fuzzy logic as one of artificial techniques to detect outliers in data set of WSN. Fuzzy logic is a logical model providing a general idea about the decision process in the analysis of the data set. The fuzzy logic suggested by [19] is essentially an approach that allows transition values to make a definition between the conventional values such as right/wrong, yes/no, high/low. The main purpose of the method is to bring a certainty to assigning a membership degree to the concepts which are hard to express or have difficult meaning. A fuzzy logic system consists of three main parts, which are fuzzification, rule base and defuzzification. Firstly, fuzzification can be defined as a transfer between a definite system and a fuzzy system and it describes a property of an object in a certain fuzzy set. The objects can belong to 'low, middle, high' property classes with membership functions, and each object is assigned to a membership degree between 0 and 1. This technique uses temporal and spatial similarity as two inputs or two membership functions to fuzzy system. These membership functions are chosen empirically and optimized using a sample input/output data. The most common membership functions include a triangle, trapezoid, Gauss curve and sigmoid. As the membership functions represent the fuzzy set, the selection of their shape and form directly affects the decision process.

Secondly, the rule base combines the membership functions from the fuzzificator with the rule handling data such as 'if, and, although, if not' which is based on the database and stored there. The If-then rules define a connecting antecedent to the consequent (i.e. input to output). These rules are given weights based on their criticality as in [19]. With this approach, measurements can be classified according to their membership degrees by adequate membership, e.g.

- If spatial link (low) and temporal similarity (low) then outlier (high)
- If spatial link (low) and temporal similarity (med) then outlier (high)
- If spatial link (high) and temporal similarity l (high) then outlier t (low)
- If spatial link (med) and temporal similarity (med) then outlier (med)

Thirdly, in the defuzzification unit, the rule results that are obtained from the rule handling unit are evaluated in the fuzzificator and turned into definite results as in [19]. Outlier is declared according to the rule results. Fig.1 represents all three stages of fuzzy logic.



**Fig. 1** Three stage of fuzzy logic

### 3.3. Outlier classification

The third step is to classify the degree of outlier value (error or event). In this step, we aim to know the source of the values labeled as outlier. There are two possible options; either this outlier value is due to an error, as a result of a low battery or network damage, or due to an event or phenomena in the surrounding environment. Our idea is based on the following observation in the result of this technique - *“Error in the sensor data are likely to be spatially unrelated while event measurements are probable to be spatially correlated”*.

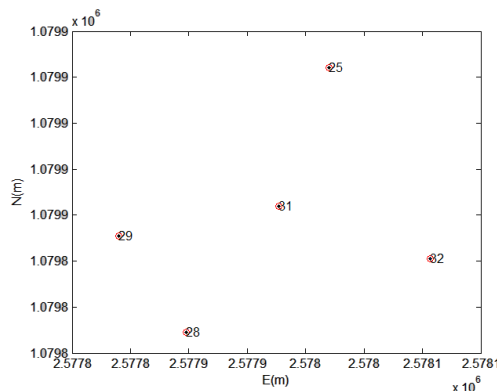
On the other hand, data instance tends to be correlated in both time and space. Hence, we employed this fact by using data from neighboring nodes to assist measuring the spatial similarity, also using time stamps between readings to assist measuring the temporal similarity. In other words, this technique detects the outlier in the previous step and if data instances are declared as outlier, it produces similar values or values larger than the outlier readings in all nodes. In addition, if those neighboring nodes readings are within the same time range, this indicates an interesting event in the physical world. Otherwise, it is likely to be an erroneous data. In our work, we assume that a sensor node (x) is considered to be a neighbor of another node (y) if x is within y’s communication range, and vice versa.

## 4. EXPERIMENTAL RESULT AND PERFORMANCE EVALUATION

In this section, we investigate the effectiveness of our proposed approach when applied on the real dataset from St.-Bernard wireless sensor network in [7]. We compare the accuracy of our algorithm with another detection method called STGOD method [15], which is based on spatial temporal correlation among neighbor nodes. We evaluate accuracy and the scalability of the proposed method against the STGOD method on a real dataset.

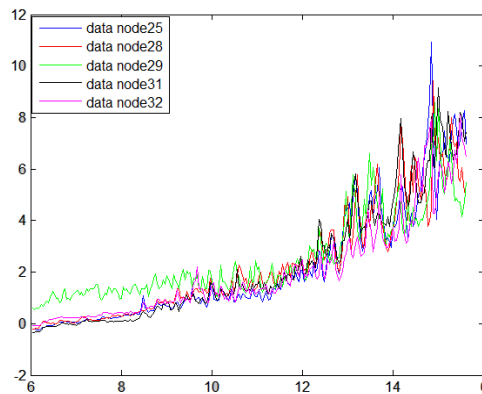
### 4.1. Study area and data description

The proposed outlier detection described in section III is applied to a realistic data set collected from 23 sensor nodes. These nodes are geographically distributed over Switzerland and Italian boarder, representing two clusters. The small cluster, situated in the Italian boarder, contains the five sensor nodes from whose data set is obtained. Fig. 2 illustrates the



**Fig. 2** A small cluster (consists of five nodes) of the Grand St Deployment and their corresponding metric coordinates (E-N).

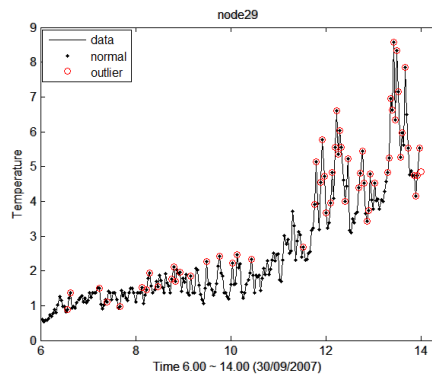
geographical distribution of these nodes over the area in which they are deployed. The collected data represent temperature as the attribute of interest. Temperature values are measured over a period 06:00–14:00 during the day (30th September, 2007). Fig. 3 depicts a plot of temperature measurements sensors for all nodes in a small cluster (node25, node28, node29, node31, node32). The measurement tolerance of the deployed sensors is about  $\pm 0.3^{\circ}\text{C}$ .



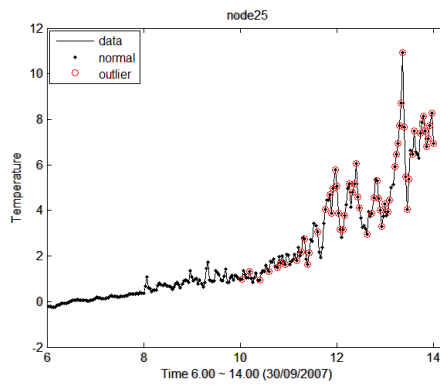
**Fig. 3** Represented data measurements of each sensor node.

#### 4.2. Results and performance evaluation

This section is devoted to evaluating the performance of the outlier detection technique proposed in section (III). Two performance metrics are considered. The first is the detection rate (DR) defined as the ratio of the correctly detected outliers to the total number of outliers in a given data set. Another performance metric of interest is the false positive alarm rate (FPR) which is defined as the ratio of normal data points incorrectly classified as outliers to the total number of normal data points. This section shows outliers in each node, detection rate, and false positive rate to each node. To evaluate performance of outlier detection needs a reference dataset. Usually, labeling techniques are utilized to label sensor measurements and classify each data point as either a normal pattern or anomalous. The choice of the labeling technique powerfully influences the evaluation of the outlier detection techniques. There are three labeling techniques used, as in [15], i.e., running average-based, Mahalanonis distance-based, and density-based, but our research used the first one which fits the data set as in [15]. In this research two software are applied, statistical model and fuzzy logic Simulink, implemented by MATLAB. As in Fig. 4 and Fig. 5, spatial temporal outliers in univariate attribute (temperature) in both node25 and node29, whose detection rate in node25 is about 92% and FPR is 10.4%, while in node29 the detection rate is 93.75% and high false positive rate is 18.33%.

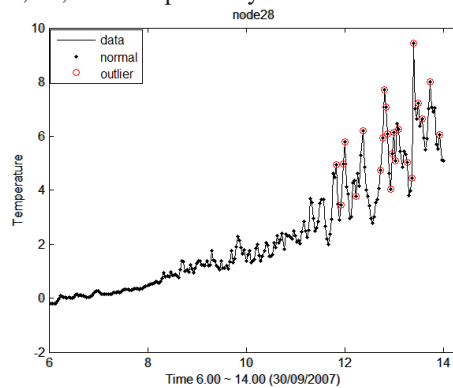


**Fig. 4** Spatial temporal outliers in node29 detected by (STODM)



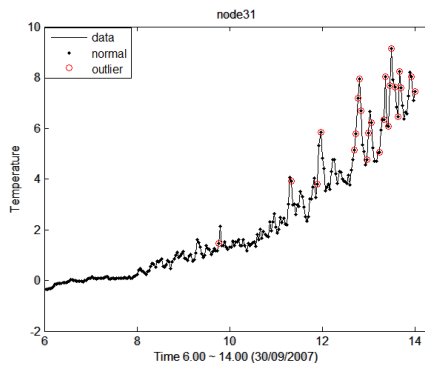
**Fig. 5** Spatial temporal outliers in node25 detected by (STODM)

While in Fig. 6, Fig. 7 and Fig. 8, node28, node31, and node32, they have high detection rate 100% and FPR 9.16, 10, 4.5% respectively in each node.

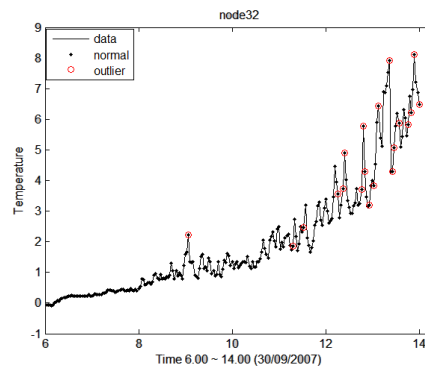


**Fig. 6** Spatial temporal outliers in node28 detected by (STODM)



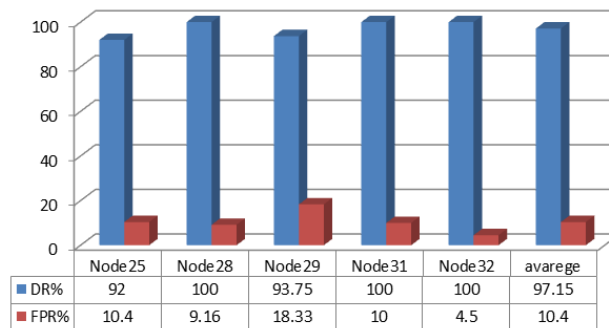


**Fig. 7** Spatial temporal outliers in node31 detected by (STODM)



**Fig. 8** Spatial temporal outliers in node32 detected by (STODM)

Fig. 9 shows the result of accuracy assessment for detected outliers by using pattern approach. The highest detection rate (100%) is at node (28, 31, 32) while the lowest detection rate (92%) is at node 25. The lowest amount of FPR is at node 32 (4.5%) while the highest rate is at node 29 (18.33%).



**Fig. 9** Accuracy of the detected outliers at different nodes

Extensive ratio on the collected data set shows that both the detection rate and FPR increase when the threshold is decreased. A fixed threshold of temporal similarity and the mean of Euclidean distance of all nodes is computed as threshold of spatial similarity that yields an average detection rate of 97.15% and FPR of 10.472%. The relative high FPR is a result of misclassifications of some normal observations, while the high detection rate achieved is a result of considering spatial temporal similarity. Table 1 shows the comparison between the proposed storms with the most frequently used data labeling technique, namely the TSOD and the STGOD technique with the detection rate and false positive alarm achieved by each algorithm. It can be observed that the proposed algorithm outperforms these techniques in terms of detection rate. Both references models are applied to the same data set as considered in our model. Another advantage of the proposed technique is that it is able to distinguish between errors and events in a given data set obtained from the sensor node. Classification of the outlier source is reported in Table 2.

**Table 1** Comparison between our approach (STODM) and STGOD model proposed of running average in [15]

Method	DR%	FPR%
STODM	97.15	10.4
TSOD	23.4	1.7
STGOD	72.34	10.94

**Table 2** Number of outliers and events detected at different nodes using STODM (our model)

Nodes	No of outlier	No of event
Node25	48	5
Node28	23	4
Node29	60	5
Node31	25	5
Node32	21	4

## 5. CONCLUSIONS

STODM algorithm proposed in this paper combines the fuzzy logic theory and distance base similarity to detect outliers and is a new try in the area of outlier detection for spatial temporal similarity. The proposed technique is able to identify normal and outlier data. Moreover, error and event are also distinguished. High detection rate is achieved compared to conventional techniques while preserving the low positive alarm rate and also reducing computational complexity because it uses Euclidian distance to calculate spatial similarity among neighboring nodes.

For future work, we plan to build an algorithm to detect outliers in multi attributes and to consider dependencies among the attributes of the sensor data as well as spatial-temporal correlations that exist among the observations of neighboring sensor nodes.

## REFERENCES

- [1] S. Kamal, R. Ramadan, F. EL-Refai, "Smart outlier detection of wireless sensor network by fuzzy logic", In Proceedings of the International Conference on Recent Advances in Computer Systems RACS-2015, Hail University, Saudi Arabia, November 2015.
- [2] Y. Zhang, M. Nirvana, H. Paul, "Outlier Detection Techniques For Wireless Sensor Networks", A Survey, University of Twente, P.O.Box 217 7500AE, Enschede, The Netherlands, 2010.
- [3] V. Chandola, A. Banerjee, A. Kumar, V, "Outlier detection: a survey", Technical Report, University of Minnesota, 2007.
- [4] V. Jha, O. Veer Singh, Y. Outlier, "Detection Techniques and Cleaning of Data for Wireless Sensor Networks", A Survey, *International Journal of Computer Science And Technology*, 2012.
- [5] X. Luo, M. Dong, Y. Huang, "On distributed fault-tolerant detection in wireless sensor networks", *IEEE Trans Computer*, Vol. 55, No. 1, pp. 58-70, 2006.
- [6] H. Konak, A. Dilaver, E. Ozturk, "The effects of observation plan and precision on the duration of outlier detection and fuzzy logic", 2005, a real network application, Survey Review, Vol. 38, 298, pp. 331-341, 2005.
- [7] Sensor Scope System. [http://sensorscope.ep.ch/index.php/Main\\_Page](http://sensorscope.ep.ch/index.php/Main_Page)
- [8] S. Subramaniam, T. Palpanas, D. Papadopoulos, V. Kalogeraki, D. Gunopulos, "Online outlier detection in sensor data using nonparametric Models", Seoul, Korea., VLDB; Young, The Technical Writer's Handbook. Mill Valley, CA: University Science, 1989, pp. 187-198M, 2006.
- [9] S. Rajasegarar, C. Leckie, M. Palaniswami, J. C. Bezdek, "Distributed anomaly detection in wireless sensor networks", UK: IEEE, ICCS, pp.12-16, 2006.
- [10] J. Branch, B. Szymanski, C. Giannella, R. Wolf, "In-Network outlier detection in wireless sensor networks", In Proceedings of IEEE ICDCS, 2006.
- [11] Rajasegarar, S., Leckie, C., Palaniswami, M. and Bezdek, J. C., "Quarter sphere based distributed anomaly detection in wireless sensor networks", Proceedings of IEEE International Conference on Communications, pp. 3864-3869, 2007.
- [12] Y. Zhang, N. Meratnia, and P.J.M. Havinga, "An online outlier detection technique for wireless sensor networks", In Proceedings of the Third IEEE European Conference on Smart Sensing and Context (EuroSSC), pp. 25-26, 2008.
- [13] Y. Zhang, N. Meratnia, and P.J.M. Havinga, "Adaptive and online one-class support vector machine-based outlier detection techniques for wireless sensor networks", In Proceedings of the IEEE 23rd International Conference on Advanced Information Networking and Applications Workshops/Symposia, pp. 990-995, 2009.
- [14] M.S. Mohamed, T. Kavitha, "Outlier detection using support vector machine in wireless sensor network real time data", *Int J Soft Comput Eng*, Vol.1, No. 2, 2011.
- [15] Y. Zhang, N.A.S. Hamm, N. Meratnia, A. Stein, M. van de Voort, P.J.M. Havinga, "Statistics-based outlier detection for wireless sensor networks", *International Journal of Geographical Information Science*, 2012.
- [16] A. Amidi, N.A.S. Hama, N. Meratnia, "Wireless Sensor Networks and Fusion of Contextual Information for Weather Outlier Detection", *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, Vol XL-1/W3, 2013.
- [17] A. Fawzy, H.M.O. Mokhtar, O. Hegazy, "Outliers detection and classification in wireless sensor networks", *Egyptian Informatics Journal*, Vol. 14, pp. 157-164, 2013.
- [18] S. Syed, M.E. Cannon, "Fuzzy logic based-map matching algorithm for vehicle navigation system", In Proceedings of the urban canyons, ION National Technical Meeting, San Diego, CA, pp. 26-28, 2004.
- [19] Y. Sisman, A. Dilaver, S. Bektas, "Outlier Detection in 3D Coordinate Transformation with Fuzzy Logic", *Acta Montanistica Slovaca Ročník 17, číslo 1*, pp. 1-8, 2012.