# Application of Modern AI Schemes for Anomaly Detection in Sensor Networks

Ashwani Kumar[1], Sumit kumar Gupta[2], Ankur shukla[3]
Electronics and Communication Engineering Department,
Bansal Institute of Engineering and Technology, Lucknow
samratashwani23@gmail.com

**Abstract: In biomedical and defense applications, ESN-based monitoring applications that make use of large data sets require outlier (anomaly) detection. Temperature, dampness, sound, strain, vibration, and other natural and actual elements are observed through a remote sensor organization. Because of the confined capacity concerning energy, memory, processing, transfer speed, the powerful idea of the organization, and the brutality of the climate, a huge assortment of sensor hubs creates a (WSN) crude sensor perceptions gathered from sensor conveys low information quality and trustworthiness. In this paper, a KNN forecast model is utilized to recognize exceptions in light of the entropy benefit of approaching sensor voltages. The calculation creation and examination depends on a constant data set delivered on 14 arrangements of MICA2 remote sensor units, with oddities entered continuously by Intel Berkeley lab volunteers involving a constant movement based interruption in the lab. To get an enormous exceptions estimations preparing dataset, a proper window size division is finished to every sensor information pair. The analysis indicates that the measurement accuracy in identifying the number of outliers is 86%. In the KNN expectation model, the calculation likewise gives an assessment of the effect of changes in distance types and number of closest neighbors.**

*Keywords—Anomaly Detection, Outliers Detection, Wireless Sensor Networks, KNN*

## 1. Introduction:

The disclosure of peculiarities, especially exceptions, is the revelation of things, occasions, or perceptions that don't adjust to an expectation in information mining. Regularly, peculiar components will cause an issue, for example, bank extortion, development deserts, clinical issues, or text based blunders. Deviations from standard schematics in the remote sensor network are alluded to as exceptions. Errors, noise, missing values, illogical data, and duplicate data are all detected by using outliers. Issues, occurrences, and malignant assaults are wellsprings of mistaken values in remote sensor organizations. An exception is a perception or mix of perceptions that digresses from the standard. An exception might be brought about by the fluctuation of the estimation. ( D. Wilson, 2014).
A straightforward outline has been attracted Figure 1 making sense of a two-layered informational index. In a wireless sensor network, the normal data regions are N1 and N2, while

points O1, O2, and O3 are outliers (S. Teng, N. Wu, and H. Zhu, 2018). These outliers fall into three broad categories: data difference, network anomaly, and node irregularity.
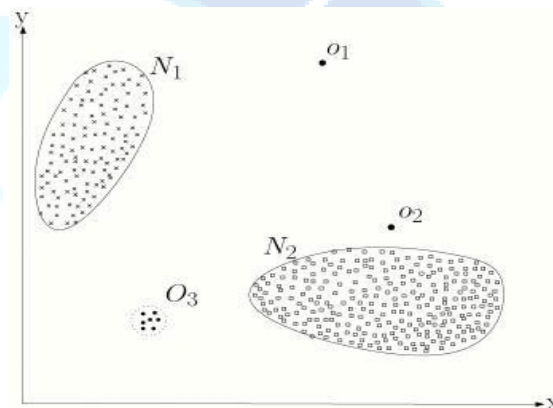


**Figure 1: Two dimensional data space clusters**

They are deployed to difficult areas, which results in these anomalies. Oddities found in a gathering of hubs are alluded to as organize peculiarities. They are the consequence of correspondence issues. Abnormalities in the network result when the sensor nodes cannot communicate with one another. DOS, sinkhole, dark opening and specific progressing and wormhole assaults are liable for these organization abnormalities. At the point when there is some anomaly in the recognized information then Information defect happens. Security infringement are additionally the justification for erroneous information. Information peculiarities are separated into three classifications, for example Transient Spatial, Spatial and Fleeting. At the point when there is change in irregularity after some time then it is known as a transient oddity. Additionally, a spatial anomaly is when an anomaly shifts to a single location in relation to a nearby node. The spatio-worldly irregularity is because of an adjustment of the worth of the information about space as well as time.

## 2. Related Work:

[4] presented Interruption Location Frameworks (IDS) for a sensor network that relies upon the association works out (e.g., number of progress and disillusionment of confirmations). In order to identify malicious offenses committed by an intruder, the framework compares occasion information and mark records.

[5] applied the revelation system in a pack based sensor network a ton of like the made structure in this composition. This sort of area system can perceive the inconsistencies that it has seen already. Regardless, this investigation is enthusiastic about recognizing peculiarities in dark circumstances, in which there are no surprising models available for the structure to learn.

[6] introduced interruption location conspires that use this model of typical traffic behavior to differentiate between unusual traffic designs. Their approaches can separate attacks that needy individual been as of late seen.

An anomaly recognition calculation based on Bayesian Belief Networks (BBN) was introduced in [8]. The structure is moreover prepared to evaluate missing characteristics in the sensor data. Both the spatial transient relationships that exist among the sensor hubs and the connection that exists between the properties of the hubs can be detected by the BBNs.

A method for revealing k-closest neighbor-based anomalies was developed by [9]: coordinates whose distance toward their k-nn outperforms a legitimate cutoff or the top n centers concerning the distance to their k-nns. Over a sliding window of its most important informative elements, each sensor keeps a histogram-like outline of the relevant data. The sink center assembles these diagrams and requests the association for any additional information expected to precisely choose the exemptions over the whole association. The usage of outlines allows less correspondence than a simple, concentrated approach. Their approach differs from ours in numerous ways. They only distinguish exceptions north of one layered information at first, and the difficulty of creating smaller, more intricate histograms will prevent any expansion beyond that. Second, they simply consider the two k-nn based exemption definitions portrayed above, while our philosophy integrates these and anything is possible from that point. Thirdly, their procedure simply applies in settings where spatial closeness is irrelevant while our system can, if fundamental, to oblige spatial area ("semi-close by" oddity ID).

Using a measure of the hidden likelihood appropriation from which the information emerges, [10] requires the sensors to maintain a tree correspondence between geography and figure anomalies. Such a check is figured by each sensor keeping an inconsistent illustration of its data insights.

[11] cultivated a framework considering a Bayesian Conviction Organization (BBN) that has been created over the IOT (and spread to each sensor). This allows each sensor to determine the likelihood of a observed tuple and, as a result, distinguish anomalies.

[12] use a wavelet-based technique for correcting colossal isolated spikes from single sensor data streams. A strong time traveling (DTW) distance-based strategy is furthermore used to perceive even more predictable time frames sensor data by differentiating the data floods of spatially close sensors expected to convey tantamount streams.

[14] focused on the idiosyncrasies in IOT, accommodating properties of irregularity acknowledgment techniques and separate the different peculiarity disclosure systems for distant sensor associations.

[15] portrayed such IOTs and the likely solutions for dealing with the recorded issues and game plan of various issues. This paper will pass on the data about the IOT and types with composing review so an individual can get more data about this emerging field.

## 3. Methodology:

The K-closest neighbors calculation (KNN) is a semi characterization method designed by Evelyn Fix and Joseph Hodges in 1951 and later extended by Thomas Covers in measurements. It is utilized in the classification and relapse of information. The source is always the K closest practical training for the data set. Contingent upon whether K-NN is utilized for classification, coming up next is the outcome:
• The result of K-NN order is a class enrollment. An item is sorted in light of a greater part vote of its neighbors, with the article being doled out to the class with the most individuals among its k nearest neighbors (K is a positive number, regularly little). • The aftereffect of K-NN relapse is the property estimation. If K = 1, the article is basically assigned to the class of that solitary closest neighbor. This number is the weighted amount of the K nearest neighbors.

K-NN is an order technique where the capability is just assessed locally. Since this technique depends on distance for arrangement, normalizing the preparation information can extraordinarily build its presentation assuming the elements address different actual units or come in various scales. ( El-Diraby Tamer E, 2020; Piryonesi S. Madeh) A good approach for both classifiers is to assign weights to the inputs of the neighbors so that the neighbors who are closer to the average contribute more than the neighbors who are farther away. A well known weighting procedure, for instance, is to provide each neighbor with a load of 1/d, where d is the distance between them ( Jaskowiak, Pablo A.; Campello, Ricardo J. G. B, 2011).

The neighbors are chosen from a gathering of things for which the class (in K-NN characterization) or cost as far as K-NN relapse is known. Albeit no express preparation stage is vital, this can be considered of as the calculation's preparation set. Uniquely, the K-NN method responds to the spatial relationship between the data.

**Algorithm:**
The preparation datset are class-named vectors in a subspace. The preparation period of the strategy comprises exclusively of saving the train tests' element vectors and enrollment capabilities.

K is a client characterized variable in the characterization stage, and a plain vector (a question or test point) is classified by giving the name that shows up most often among the K marked information nearest to that keypoint.

Time series distances are frequently measured in Euclidean units. Another measurement, like the covers metric, can be utilized for unmitigated information, like text arrangement (or Hamming distance). K-NN has been utilized as a measurement with connections, for example, Pearson and Spearman in the setting of quality articulation microarray information. ( Pablo A. Jaskowiak, 2011) When the distance measure is learned with particular calculations like Huge Room for error Closest Neighbor or Nearby parts investigation, the grouping execution of K-NN can frequently be significantly moved along.

The fundamental categorization of "majority voting" suffers from a flaw when the classifier is distorted. That is, due to their extraordinary amount, models from a more continuous class will generally win the figure of another model. ( Brian S. Everitt, Sabine Landau, and MorvenLeese. 2011) One option is to weight the categorization by taking into account how far apart the test location is from each of its k closest neighbors. Every one of the k closest focuses' class (or worth, in relapse issues) is figured as the aggregate corresponding towards the point relative between that point and the test point. Straightforwardness in gathered information is one more strategy to defeat slant. In a self-sorting out map (SOM), for instance, every hub addresses (or fills in as the focal point) of a group of comparative times, no matter what their thickness in the preparation information. The SOM can then be dealt with utilizing K-NN.

**Determination of boundaries:** The ideal decision of not entirely settled by the information; by and large, more prominent upsides of K lower the impact of commotion on recognizable proof, however make class limits less evident (Nigsch, Florian; Drinking spree, Andreas; van Buuren, Bernd, 2006). Different heuristic systems can be utilized to choose a decent K. (see hyper boundary streamlining). The closest neighbor calculation is utilized when the class is projected to be the class of both the nearest base classifier (for example at the point when K = 1).

The precision of the K-NN strategy can be significantly hurt in the event that there are loud or immaterial highlights present, or on the other hand in the event that the trademark scales are not corresponding to their significance. To further develop classification, a ton of review has performed into choosing or

scaling qualities. The applicatin of transformative calculations to boost include scaling is an especially normal technique. ( Peter Lobby, Byeong U. Park, and Richard J. Samworth, 2008) Another normal strategy is proportional highlights utilizing shared data between the preparation information and the instructional courses.

While managing parallel (two-class) text arrangement, it's ideal to make K an odd number to stay away from tying votes. The bootstrap approach is a famous approach to deciding the tentatively ideal K in this present circumstance (Stone, Charles J, 1977).
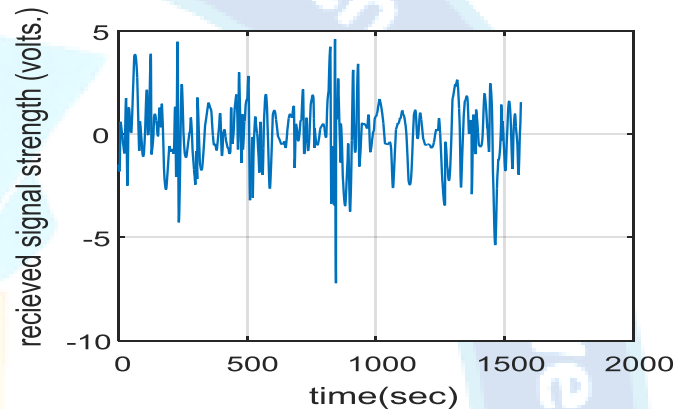
**4. Result and Discussion:**
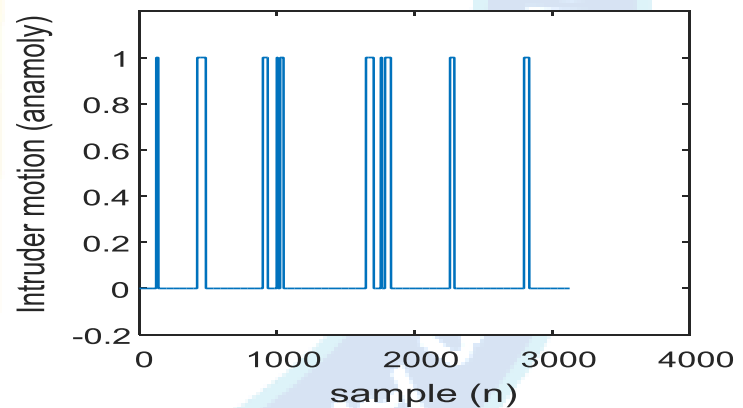


**Figure 2: Received signal strength vs time**



**Figure 3: Intruder motion (anamoly) vs sample**

Figure 6 shows the plot up between greatest preparation and testing exactness acquired during the endeavor 1 of KNN based location. The red line shows the preparation precision result at various number of neighbors (x pivot). The quantity of neighbors are shifted from 2 to 10.It has been seen that most noteworthy preparation exactness in neighbor wise is at 2 closest neighbors having esteem 83.39.The blue line addresses the preparation precision at the distance wise (D1 to D10).In this case at the D2 the preparation exactness lies for 2 number of closest neighbor.
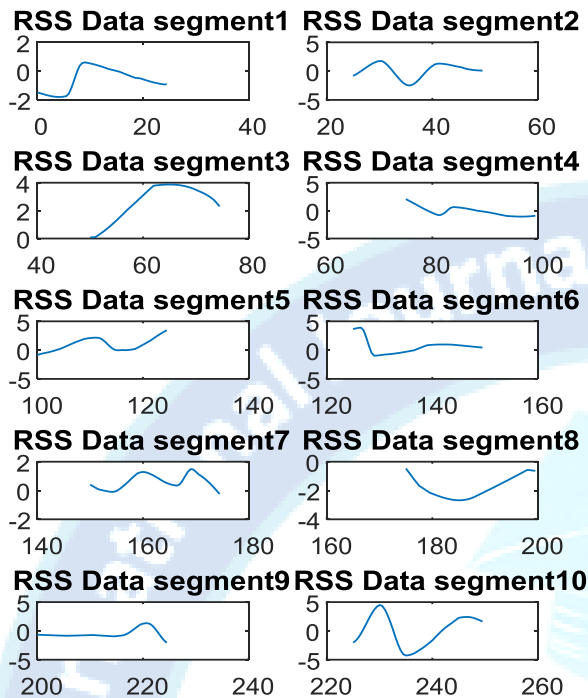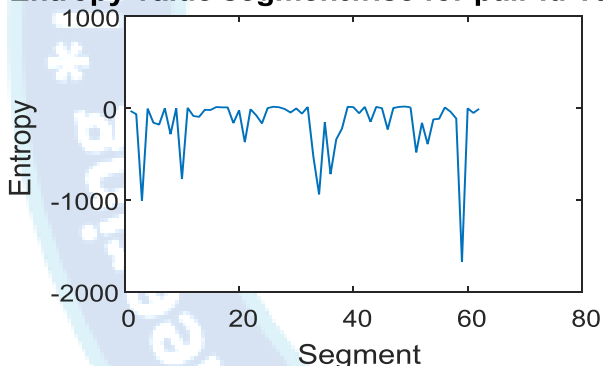
Figure 4: RSS data segments of received data at length of 50 samples.



**Figure 6: Maximum accuracy**

The testing precision in rate is shown by line green for sign number of closest neighbor and dark for distance type. It has been seen that the most elevated testing precision lies at 76.89% for distance type D7 for 9 number of neighbors as found in neighbor wise plot.

**5. Conclusion:**
This paper introduced the approach for detection of outlier using KNN approach. It describes the use of sensor data records affected by outlier due to intrusion. It explains the steps included in data acquisition and segmentation. After pre-processing the data is passed through the process of feature extraction to get the segment wise entropy values. Finally it describes the process of developing model K nearest neighbour for outlier detection using the entropy feature.

**References:**
[1] Chandola, V.; Banerjee, A.; Kumar, V. Anomaly detection: A survey. ACM Comput. Surv. 41, 15, 2009.
[2] Raja Jurdak, X. Rosalind Wang, Oliver Obst, and Philip Valencia,"Wireless Sensor Network Anomalies: Diagnosis and Detection Strategies", 2011.
[3] Y. Zhang, N. Meratnia, and P. J. M. Havinga, "Outlier Detection Techniques for Wireless Sensor Network" A Survey,Technical Report, University of Twente, 2008.
[4] Techateerawat, P. and Jennings, "A Energy efficiency of intrusion detection systems in wireless sensor networks", In Proceedings of the 2006 IEEE/WIC/ACM international conference on Web Intelligence and Intelligent Agent Technology (WI- IATW), pages 227–230, Washington, DC, USA. IEEE Computer Society, 2006.
[5] Hai, T. H., Khan, F. and nam Huh, E. "Hybrid intrusion detection system for wireless sensor networks", Computational Science and Its Applications ICCSA, 4706:383–396, 2007.
[6] Onat and Miri, Onat, I. and Miri, "An intrusion detection system for wireless sensor networks", In IEEE International

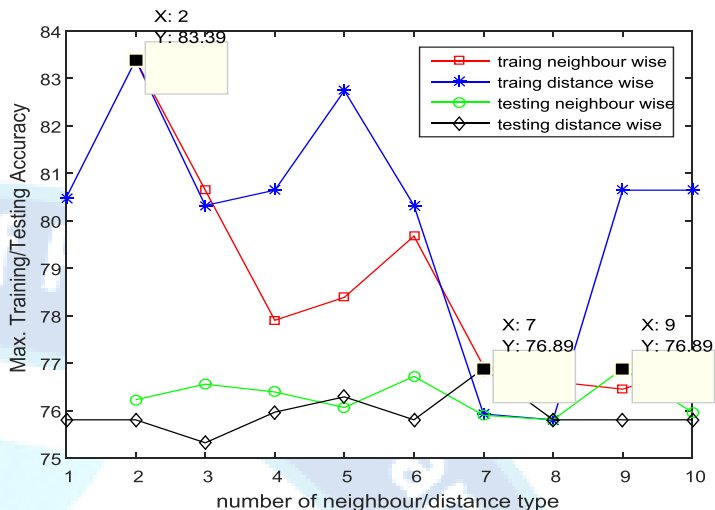**Figure 5: a) Entropy value vs segment of pairid 167, b)Anamoly level segmentwise vs segment of pairid 167.**

Conference on Wireless And Mobile Computing, Networking and Communications, Los Alamitos, IEEE Computer Society Press, 2005.

[7] Banerjee, Banerjee S., Grosan C., and Abraham, "Ideas: intrusion detection based On emotional ants for sensors", In The 5th International Conference on Intelligent Systems Design and Applications (ISDA), pages 344–349, Wroclaw, Poland, 2005.

[8] Janakiram, D., Reddy, V., and Kumar A., "Outlier detection in wireless sensor Networks using Bayesian belief networks". In First International, Conference on Communication System Software and Middleware, pages 1–6, 2006.

[9]. Patcha, A.; Park, J.-M. An overview of anomaly detection techniques: Existing solutions and latest technological trends. Comput. Netw. 2007, 51, 3448–3470.

[10]. Snyder, D. Online Intrusion Detection Using Sequences of System Calls. Master's Thesis, Department of Computer Science, Florida State University, Tallahassee, FL, USA, 2001.

[11]. Markou, M.; Singh, S. Novelty detection: A review—Part 1: Statistical approaches. Signal Process. 2003, 83, 2481–2497.

[12]. Markou, M.; Singh, S. Novelty detection: A review—Part 2: Neural network based approaches. Signal Process. 2003, 83, 2499–2521.

[13]. Hodge, V.; Austin, J. A Survey of Outlier Detection Methodologies. Artif. Intell. Rev. 2004, 22, 85–126.

[14]. Goldstein, M.; Uchida, S. A Comparative Evaluation of Unsupervised Anomaly Detection Algorithms for Multivariate Data. PLoS ONE 2016, 11, e0152173.

[15]. Wang, H.; Bah,M.J.; Hammad,M. Progress in Outlier Detection Techniques: A. Survey. IEEE Access 2019, 7, 107964–108000.

[16]. Tellis, V.M.; D'souza, D.J. Detecting Anomalies in Data Stream Using Efficient Techniques: A Review. In Proceedings of the 2018 International Conference on Control, Power, Communication and Computing Technologies (ICCPCCT), Kannur, India, 23–24 March 2018.

[17]. Park, C.H. Outlier and anomaly pattern detection on data streams. J. Supercomput. 2019, 75, 6118–6128.

[18]. Pimentel, M.A.; Clifton, D.A.; Clifton, L.; Tarassenko, L. A review of novelty detection. Signal Process. 2014, 99, 215–249.

[19]. Chauhan, P.; Shukla, M. A review on outlier detection techniques on data stream by using different approaches of K-Means algorithm. In Proceedings of the 2015 International Conference on Advances in Computer Engineering and Applications, Ghaziabad, India, 19–20 March 2015.

[20]. Salehi, M.; Rashidi, L. A Survey on Anomaly detection in Evolving Data. ACM Sigkdd Explor. Newsl. 2018, 20, 13–23.