

Outlier Detection in Wifi Indoor Localization and Internet of Things (IOT) Using Machine Learning

Nirmal Kumar M

Computer Science and Engineering &

Sethu Institute Of Technology

Abstract - In Internet of Things (IOT) millions of devices are connected for providing smart services. In indoor localization environment, that is one of the most concerning topic of smart cities, internet of things and wireless sensor network. Wifi indoor localization environment by analyzing RSSs using the combination of supervised, unsupervised and ensemble machine learning methods. The input dataset is in the format of .csv file. In data pre-processing, data has been shifted from text file to .csv because .csv file provide well representation of data in the shape of rows and columns. In the .csv file heading for seven columns containing RSS value. Was given as rssi to rssi7 and for heading of class label 'y' has been used. In this paper, KNN and Isolation Forest are used as supervised learning. The method is proved to be effective with high accuracy, precision, recall, F-score.

Key Words: Internet of Things, Localization, outliers, outliers detection.

1. INTRODUCTION

Indoor localization is a process through which nodes of a network obtain their location in indoor environments such as smart homes, using smart devices. The term outlier, also known as anomaly is originally taken from the field of statistics. Outliers can be raised because of human error, machine error, mechanical fault and changes in the behaviour of system or may be due to natural deviance in the environment. While in the case of Wi-Fi indoor localization environments, outliers are irregular and distinct RSS values caused in Wi-Fi indoor localization environment. To detect the outliers in the data of RSSs this research proposed to use unsupervised learning method such as isolation forest (I Forest) to initially classify the data into two classes as normal and abnormal. After that machine learning classifiers support vector machine (SVM), K-Nearest neighbour (KNN). The dataset was collected by observing signal strengths of seven Wi-Fi routers using smartphone and saved in a text file to perform experimentation of indoor localization. The process of finding the spatial location of nodes in a wireless network has been called localization, positioning, geo location, and self-organizing in the literature. The term localization is the most popular and so is used here. In a wireless network, we

can classify the nodes into three categories, anchor, un-localized, and localized. The first group of nodes know their position or coordinates and are called *anchors*. Nodes in the second group do not know their position and are called un-localized. The third group contains those nodes which were in the second group but subsequently had their positions estimated, and thus are called localized. The location information of nodes in a wireless network can be used for many useful purposes such as tracking mobile nodes, determining the coverage area, load and traffic management, node lifetime control, cluster formation, and routing enhancement.

A. Accuracy of Localization

RSSs of a Wi-Fi network are actual signal strengths measured at receiver's device in decibel-milliWatts (dBm) or milliWatts (mW). RSS based techniques for indoor localization, such as Wi-Fi and bluetooth low energy (BLE) device are easy to deploy and cost efficient. Their signal strengths can be affected by the external forces such as reflection, refraction, diffraction, interference etc. due to the inherent nature of signals and from physical obstructions. The reference node in the wireless environment will receive different phase delayed and attenuated versions of the same signal. Almost all the techniques used for localization such as time difference of arrival (TDoA), time of flight (ToF), angle of arrival (AOA), including RSS depends upon these signals from the user's device or a reference node. A phase delayed and attenuated signal affected by reflection, refraction, diffraction, interference etc. causes irregularity in the RSS at reference node. This irregularity of RSS values affects the accuracy of localization and positioning in indoor environment.

B. Radio Environment

Wi-Fi based indoor networks use radio signals for communication. Radio signals are greatly affected by the physical obstacles such as walls, roof, doors and even from the human beings. Electric and magnetic field created by other electronic devices that are being used in indoor environment also effect radio signals. Radio signal strengths received at reference node also depend upon the location of node. These discussed factors about radio environment should essentially be taken into account for an accurate localization in indoor environment.

C. Bad Positioning

Bad positioning is actually incorrect position or location in indoor environment. For example, a smart device's user is trying to reach in the admin office of a department that is located on the second floor of a large building. A good indoor location system will show him exact position but a bad location system may show his location on first floor, third floor or may be an unexpected location. Hence, it's good to say that without positioning is better than bad positioning. Bad positioning in Wi-Fi indoor localization system can also be due to the irregular and anomalous RSSs caused by the external effects and physical obstructions of indoor environment.

Therefore, this research presents an outlier detection technique "iF_Ensemble" for Wi-Fi based indoor localization environment using the combination of supervised, unsupervised and ensemble learning methods. To detect the outliers in the data of RSSs this research proposed to use unsupervised learning method such as isolation forest (iForest) to initially classify the data into two classes as normal and abnormal. After that machine learning classifiers support vector machine (SVM), K-nearest neighbor (KNN), random forest (RF) and Elliptic envelop has been applied separately on the data. However, their results were not satisfactory. Therefore, stacking an ensemble learning method has been proposed for the improved separation of outlier and normal RSSs.

The rest of the paper is organizing as follows. Section II addresses related work of outlier detection. Section III describes material and method used for outlier detection. Section IV describes the results and discussion. Finally, Section V concludes this paper.

2. RELATED WORK

Data analysts have always been concerned with the detection and correction of outliers of data because of security threats and unusual events. Outlier detection is an imperative area in studies and research of data mining and have various applications such as weather prediction, fraud discovery and interruption detection for a networks [6]. Detection of outlier has long been a field of study and research. An outlier detection technique called localization anomaly detection (LAD) was proposed in [7]. Their outlier detection scheme

tries to detect the outliers and perform negotiation resilient localization without eliminating the anomalous node from the network. An anomaly detection technique to deal with the anomalous sensor's data has been presented in [8]. Authors discussed that low cost RSSs are mostly used for localization indoor environment but these RSSs are greatly affected by the environment changes. They presented their technique for outlier detection purpose to deal with anomalous RSS data in indoor localization so that they can obtain reliable measures of RSS data for localization purpose.

An algorithm based on robust RSS for a secure localization has been presented in [9]. The proposed algorithm uses robust localization and time synchronization primitive, which correctly joined and allows attack's detection in the localization within an accurate attacker's model. A technique for outlier detection measurements that depends upon graph rigidity has been proposed in [10]. In the proposed approach authors examine that how much knowledge is required for the identification of outlier measurements. They propose the idea of confirmable edges by using the rigidity theory, and also develop the terms for an edge to be confirmable. A range based algorithm for outlier detection called localization density based outlier detection (LD-BOD) has been proposed that deals with the premise of density based anomaly detection strategies from the field of data mining [1]. To select a paramount candidate point, density based outlier detection uses the distance to K-nearest neighbor that is an average used for the estimation of location of unidentified node.

An Outlier detection approach using Bayesian belief networks (BBNs) for wireless sensor networks (WSNs) has been proposed in [11]. In their proposed approach authors discussed that in the user's perspectives; data consistency is an imperative problem in the context of WSNs steamed data. The consistency of data is affected by the severe condition of environment, interference in radio signals and bad quality sensor's usage. So, the data produced by the sensors can be incorrect and inconsistent that is main cause of missing values and outliers in the data. A clustering based outlier detection algorithm for WSNs called outlier detection based on clustering (ODC) has been presented in [12]. According to authors their proposed algorithm (ODC) will catches hidden outliers by the detection of labels of cluster in the slots of time.

An anomaly detection approach using improved iForest has been proposed in [13]. The proposed approach used simulated annealing method for the optimization of iForest. This simulated annealing method improved the speculation ability of iForest, eliminates terminated isolation tree (iTree) to minimize the computational issues. It improved speed of detection, stability and enhanced the performance of data outlier detection. An IoT architecture for the detection of outliers in the dataset of a forest environment has been proposed in [14] by applying four models of statistics such as gradient boosting machine (GBM), RF, classification and regression trees (CARD) and linear discriminant analysis (LDA).

An effective multi-dimensional process for user behavior anomaly detection (UBAD) has been

developed using test of multi-dimensional statistics for anomaly detection [15]. They discussed that user behavior outlier detection is an imperative phase of user behavior analysis that can be used for the detection of events of an anomalous user from network traffic data and event. A methodology for the detection of abnormal actions of devices by examining their data of events has been proposed in [5]. Authors discussed that devices are associated with the internet and they are sending their log information to the server.

They mainly focus on the analysis of this data for automatically detecting infrequent patterns. They apply three steps i.e., feature or attribute generation, attribute aggregation and data analysis to process the log data. According to best of our knowledge, in all above discussed previous studies of outlier detection in indoor localization and Internet of Things, machine learning methods have not been used along with ensemble learning as proposed in our research.

3. MATERIAL AND METHODS

A. The Data Set

The test dataset that has been used in the proposed iF_Ensemble technique for anomaly detection purpose has been taken from openly available database repository provided by the University of California, Irvine (UCI)¹. The dataset was collected by observing signal strengths of seven Wi-Fi routers using smartphone and saved in a text file to perform experimentation of indoor localization. The dataset contains 2000 entries (rows) with eight features (columns). Seven columns of the dataset contain the values of RSSs of seven Wi-Fi routers. The eighth column is class label as 1 to 4 that corresponds to different locations of an office such as conference room, kitchen or the indoor sports room etc.

B. Methods Used for Outlier Detection and Indoor User's Localization

Techniques used in the proposed approach include following machine learning methods.

1) Data Preprocessing

Raw data is mostly incomplete, inconsistent and may not be in a suitable format to perform an action. Data preprocessing is performed to clean, transform or normalize this data to bring it in the required and useful format.

In data preprocessing, data has been shifted from text file to.csv because .csv file provide well representation of data in the shape of rows and columns. In the .csv file headings for seven columns containing RSS values was given as rssi1 to rssi7 and for heading of class label 'y' has been used, so that it can be easily used for coming stages. After that .csv file was imported in a python program using python libraries such as pandas and required columns containing the values of RSS has been selected using columns function. At the end transformation of data into two dimensions (2d) has been made

using decomposition library, principal component analysis (PCA) and its transform function in python.

2) Supervised Learning

Supervised learning is training of machine using some data that has already been labeled correctly. In this learning style corresponding proper outputs are provided along with the instances [16]. Some of the important supervised learning techniques that are proposed for outlier detection purpose in this research are:

- SVM with the parameter setup as $kernel='rbf'$ and $gamma = 0.0001$. $tridSearch()$ has also been used for the optimization of the parameters but the above given parameter setup is best for the underlying dataset.
- RF with parameter as $random_state = 0$, $n_estimators = 100$ and $maxdepth = 2$.
- K-nearest neighbors with $n_neighbors = 5$.

3) Unsupervised Learning

Unsupervised Learning is the training of a machine using information that is neither labeled nor classified. In this learning style corresponding proper outputs are not provided along with the instances [16]. As in unsupervised learning there is no need of human expert to manually label the data so it is extensively appropriate than supervised learning. Unsupervised learning method used in this research includes iForest with the parameter setup such as $n_estimators = 5$ and $contamination='auto'$. Local outlier factor (LOF) also has been used to compare with proposed iForest with $n_neighbors = 5$.

4) Ensemble Learning

Ensemble learning is a machine learning technique used to combine multiple weak learners into one but more effective model to get the better results [17]. It is a technique that uses multiple machine learning models instead of single model to solve a particular problem. The three terms used to describe the ensemble learning are Bagging, Boosting and Stacking. In this research we prefer to use stacking that is used for combining different classifiers to improve their predictive force. It is used to combine and learn numerous different weak learners by training a meta-model for the final predictions results. This meta-model based on the predictions results of weak learner models for the final output.

C. Evaluation Methods

For the evaluation and accuracy checking of outlier detection methods as well as classification methods of the proposed approach, following method has been used:

1) Confusion Matrix

Confusion matrix is used to measure the performance of classifier in machine learning on a test set of data that's true values are already known. The proposed approach of outlier detection has been implemented in python programming and python programming provide a library named metrics and the classes in it $confusion_matrix()$ and

classification_report(). These python classes take the actual class labels of the underlying problem and the predicted result of classifier as an input and returns in a form of matrix. Results of confusion matrix are used to calculate other performance metrics such as accuracy, recall, precision, F-score and specificity.

2) Receiver Operating Characteristic (ROC) curve

ROC curve is a graph that is used for the evaluation of classifier's performance in machine learning. It plots the true positive rate (TPR) on y-axis versus false positive rate (FPR) on x-axis at different classification thresholds. It is a probabilistic curve and specifies an area under curve (AUC) that represents the degree of separability. ROC (AUC) is mostly used when classes of the underlying problem are balanced.

3) Precision-Recall curve

A Precision-Recall curve is a graph used for the evaluation of classifier's performance in machine learning using precision and recall scores. Precision is also known as Positive Predictive values while Recall is also known as sensitivity or TPR. In the plot of Precision-Recall, Recall or TPR is placed on x-axis while precision is placed on y-axis. It is mostly used when classes of underlying data are unbalanced.

4. IMPLEMENTATIONS

The proposed approach of outlier detection iF_Ensemble has been implemented in Python programming language using tools Jet Brains PyCharm and Jupyter notebook. Python is most famous programming language that is used in the field of data science because it provides numerous libraries, built in classes and function for data processing and visualization. Following steps have been performed in iF_Ensemble using python program for detection of outliers from the RSS data of indoor localization environment.

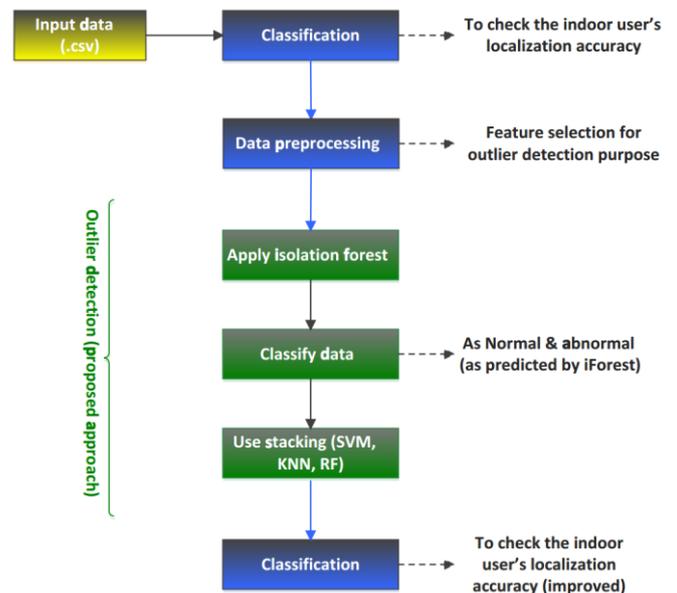
Input data is taken as a .csv file. To check the localization accuracy, user's localization is performed using classification methods KNN, SVM and Naïve Bayes according to the class labels (1 to 4) used for localization. Data preprocessing is performed to select the required features for outlier detection. These features are seven columns containing RSSs values. iForest an unsupervised machine learning method is then applied on the selected features

Data is classified as normal with class label 1 and abnormal with class label 1 as predicted by the iForest. Stacking, an ensemble learning technique that is used to improve the performance of classifiers is then

applied. SVM, KNN and RF have been used as base learner and RF as meta learner for stacking. Elimination of the data observation predicted by stacking as outliers is performed. To check the improvement in localization accuracy after applying the outlier detection technique, user's localization is performed using classification methods KNN, SVM and Naïve Bayes.

5. RESULTS AND DISCUSSION

In this section calculation of results and evaluation of the proposed approach has been made using the test dataset that was collected in indoor environment. A python code was implemented for this purpose and it was executed tentimes/iterations. For each iteration; precision, recall, F1-score and accuracy was computed from the results of



. Fig 1 .Semantic diagram of proposed outlier detection approach iF_Ensemble.

Table 1. Results of user's localization before applying the outlier detection approach.

Model	Precision	Recall	F1-score	Accuracy
SVM	96.41%	96.33%	96.25%	96.4%
KNN	96.45%	96.22%	96.22%	96.6%

Table 2. Results of unsupervised learning methods.

A. Indoor User’s Localization Accuracies before Applying the Proposed Outlier Detection Approach

Classification methods KNN [18], SVM [19] and Naïve Bayes [20] are performed on the indoor dataset to observe the localization accuracy in Wi-Fi connected indoor environment. The results of this step were later used to measure the difference in localization accuracy obtained after applying the proposed approach of outlier detection. Table 1 shows the accuracy results of user’s localization in Wi-Fi indoor localization environment before applying the proposed outlier detection approach.

B.Results of Proposed Outlier Detection Technique:iF_Ensemble

In iF_Ensemble, initially seven columns containing RSS values of Wi-Fi routers has been selected. As in case of raw data there was no clue that which portion of data is normal and which is abnormal (outliers). Therefore, iForest and LOF, two unsuper-vised learning methods, were applied on the selected features. As it is evident from Table 2 that accuracy of iForest is better than LOF thus we proposed to use iForest for separating normal data from abnormal in our proposed outlier detection technique.

In iForest result is predicted in the form of 1, 1 where 1 represents normal and 1 was used for abnormal data. Under-lying dataset was assigned to a dataFrame in python program.

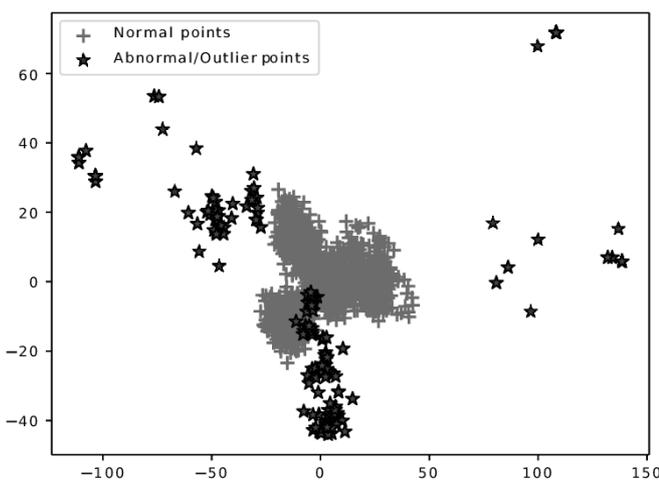


Fig 2 . Scatter plot of normal and abnormal classified data using iForest.

Model	Accuracy
iForest	90.75%
Local outlier factor	88.45%

Table 3 Result of outlier detection

Model	Precision	Recall	F1-score	Accuracy
SVM	96.18%	94.57%	95.64%	96%
Elliptic envelope	93%	91.93%	92.32%	92%
RF	94.79%	94.57%	93.73%	94.6%
KNN	96.89%	96.86%	96.85%	96.71%
Stacking	97.62%	97.67%	97.62%	97.8%

and labels 1 and 1 (predictions of iForest) were concatenated with the dataFrame against each entry (row) as eighth column. This eighth column is used as class label so that we can apply the classification methods for outlier detection. The results of iForest have been visualized in Fig. 2 using a two dimensional scatter plot of normal and outlier data.To increase the accuracy level of proposed outlier detection method stacking is used. The main reason for using stacking is to enhance the accuracy of abnormal data identification with the help of classification methods.Stacking uses base learner for initial predications and then a meta learner for final predication. Therefore, in this step of pro-posed approach KNN, SVM and RF has been used as base learn-ers for stacking, while RF has been used as meta learner for the final predication to achieve the better results of outlier detection. These machine learning classifiers such as SVM, KNN, RF in-cluding elliptic envelope has also been applied separately (one by one) on the dataset to compare their results with proposed stacking.Allused models including stacking has been trained with 80 percent of the data observations and then whole dataset was tested using trained models instead of 20 percent of the remain-ing test data because our task was to differentiate outlier and normal data from the whole dataset. It is evident from Table 3 that stacking method showed highest accuracy of 97.8%.Fig. 3 shows the accuracy of the algorithms used for outlier detection. Stacking algorithm achieves highest accuracy among all other algorithms.We did not only evaluate the results of outlier detection using metrics such as precision, recall, F1-score and accuracy but also evaluates the results by plotting Receiver Operating Character.

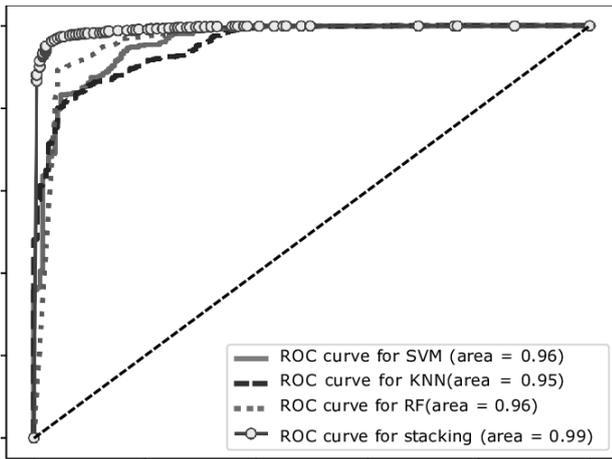


Fig 3 ROC-AUC curve of outlier detection using SVM, ELPEN, RF and Stacking technique.

Istic curve that is shown in the Fig. 4. The ROC graph shows the ROC curves and AUC of four different algorithms applied on the dataset that was classified as normal and abnormal by iForest. ROC curve shows the measurement of sensitivity con-sidering ROC curve by plotting a function with (1-specificity) and the point on ROC expresses as sensitivity/specificity pairs those are comparing with a threshold [21]. The AUC and results will be better and high when a ROC curve passes from (0, 1) coordinate or ROC curve passes from near the upper left corner.ROC graph plotted in the Fig. 4 shows 95 percent AUC for of elliptic envelope and 96 percent for both SVM and RF. It also shows ROC-AUC of proposed approach that is stacking with almost 99 percent AUC that is better than others.

The results of outlier detection have also been evaluated using precision-recall curve. Precision-recall curve is plotted inFig. 5 that includes the precision-recall curves of SVM, elliptic envelope and RF. It also shows precision-recall curve of proposed approach that is stacking. Precision-recall curve also

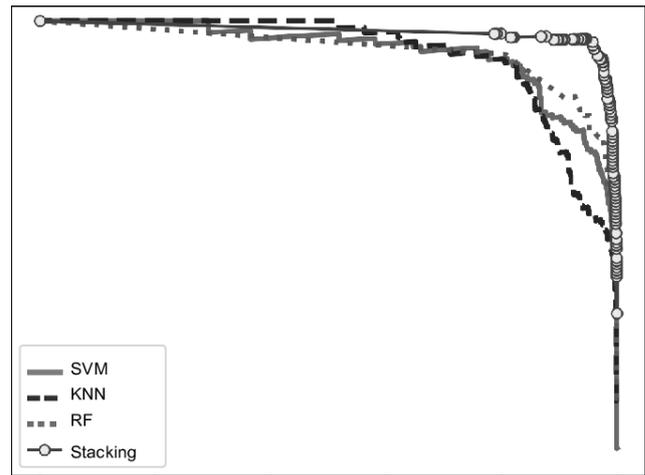


Fig 4 Precision-recall curve of outlier detection using SVM, ELPEN, RF, and stacking technique.

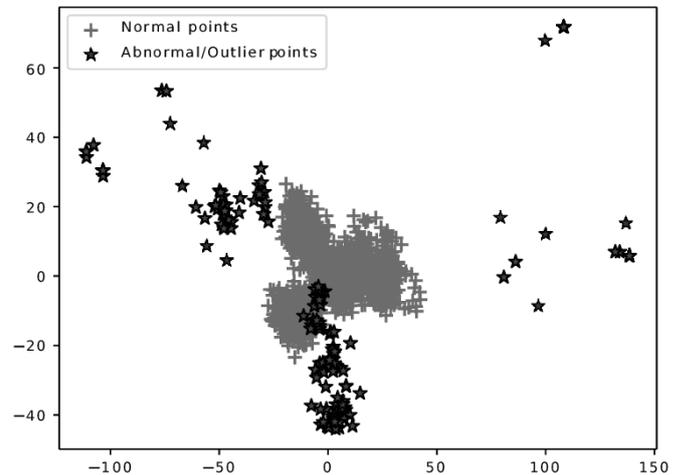


Fig. 6. Scatter plot of normal and abnormal classified data using iF_Ensemble, the proposed stacking based outlier detection approach

indicates the better results of stacking.

After separating the normal and outlier data in the proposed approach outlier and normal data has been transformed into two dimension using PCA and transform related functions in python. Then these transformed values have been plotted in a scatter plot to visualize the result of proposed outlier detection approach as shown in Fig. 6. Figure shows that normal data points that were declared abnormal by iForest (Fig. 2) have been correctly identified as normal by applying stacking. Thus the proposed approach "iF_Ensemble" detects outliers more efficiently.

C. Indoor User’s Localization Accuracies after Applying the Proposed Outlier Detection Approach iF_Ensemble

After applying the iF_Ensemble approach, outlier observation has been eliminated by skipping the observation contain-

Table 4. Results of user's localization after applying the outlier detection approach.

Model	Precision	Recall	F1-score	Accuracy
SVM	98%	98%	98%	98.30%
KNN	98.25%	98.18%	98.20%	98.50%

Table 5. Accuracy comparison of user's indoor localization.

Model	Before	After
SVM	96.4%	98.4%
KNN	96.6%	98.3%

Table 6. Comparative results based on detection rate/recall.

Model	Detection rate/recall
SVM	93%
Neural network	92%
Decision tree	87%
iF_Ensemble(Proposed)	97.67%

classifiers that were previously used for user’s localization in in-door environment, i.e., SVM, KNN, and Naïve Bayes have been used. Table 4 shows the results of user’s indoor localization after applying outlier detection approach. Comparison of accuracies of indoor user’s localization using classification methods such as SVM, KNN and Naïve Bayes before and after applying the proposed outlier detection approach is shown in the Table 5. It clearly indicates that accuracy of indoor localization improved after applying the proposed iF_Ensemble approach.

D. Comparative Analysis

To confirm the accuracy and performance of the iF_Ensemble it has been compared with the results of classification based methods for outlier detection in [22]. In which authors provides a complete and structured overview of classification based outlier detection techniques used in Wireless Sensor Networks. Comparative results based on detection rate (DR) or recall of the proposed approach and other classification based algorithms are shown in Table 6 where better results of our proposed approach can be observed.

6. CONCLUSIONS

In indoor localization and IoT location information plays imperative role for accurate localization. But the RSSs value in Wi-Fi indoor environment are affected by the limitation of radio signals that causes irregularity in RSSs. These irregular and anomalous signals could not correctly identify a node in a local network. Therefore, this research has developed an outlier detection technique named as "If Ensemble" for a Wi-Fi indoor localization environment by analyzing RSSs using machine learning methods. The results evaluation of proposed approach using precision, recall, F-score and accuracy showed better results with 97.62 percent of precision, recall, and F-score and with 97.8 percent accuracy. The AUC of ROC and precision-recall curves plotted for the evaluation also confirm the better results of the proposed approach with almost 99 percent of AUC. In this research the accuracy results of user’s indoor localization has also been calculated before and after outlier detection and elimination of outlier. It is observed that indoor user’s localization accuracy improved by 2 percent. So, it is concluded that this can be useful and reliable approach for outlier detection in indoor localization and IoT. As general, there are little techniques for outlier detection in localization environments and wireless sensor networks. The applicability of artificial intelligence (AI) techniques for outlier detection in localization and WSNs can be investigated in the future.

REFERENCES

- K. K. Almuzaini and A. Gulliver, "Range-based localization in wireless networks using density-based outlier detection," *Wireless Sensor Network*, vol. 2, no. 11, p. 807, Nov. 2010.
- H. Aly, A. Basalamah, and M. Youssef, "Accurate and energy-efficient GPS-less outdoor localization," *ACM Trans. Spatial Algorithms Syst.*, vol. 3, no. 2, pp. 1–31, July 2017.
- D. Janakiram, V. A. Reddy, and A. P. Kumar, "Outlier detection in wireless sensor networks using Bayesian belief networks," in *Proc. COMSWARE*, 2006.
- Y. Zhang, N. Meratnia, and P. Havinga, "Outlier detection techniques for wireless sensor networks: A survey," *IEEE Commun. Surveys & Tuts.*, vol. 12, no. 2, pp. 159–170, 2010.
- K. Niu, F. Zhao, and X. Qiao, "An outlier detection algorithm in wireless sensor network based on clustering," in *Proc. IEEE ICCT*, 2013.
- D. Janakiram, V. A. Reddy, and A. P. Kumar, "Outlier detection in wireless sensor networks using Bayesian belief networks," in *Proc. COMSWARE*, 2006.
- K. Niu, F. Zhao, and X. Qiao, "An outlier detection algorithm in wireless sensor network based on clustering," in *Proc. IEEE ICCT*, 2013.
- X. Wei, L. Wang and J. Wan, "A New Localization Technique Based on Network TDOA Information," *Proceedings of IEEE International Conference on ITS Telecommunications*, Chengdu, China, June 2006, pp. 127-130.
- Y. Zhang, N. Meratnia, and P. Havinga, "Outlier detection techniques for wireless sensor networks: a survey," *IEEE Communication Surveys & Tutorials*, vol. 12, no. 2, 2010.
- X. Wei, L. Wang and J. Wan, "A New Localization Technique Based on Network TDOA Information," *Proceedings of IEEE International Conference on ITS Telecommunications*, Chengdu, China, June 2006, pp. 127-130.
- M. Ahmed, A. N. Mahmood, and J. Hu, "A survey of network anomaly detection techniques," *J. Network Computer Applicat.*, vol. 60, pp. 19–31, 2016.
- A. Ayadi et al., "Performance of outlier detection techniques based classification in wireless sensor networks," in *IEEE IWCMC*, 2017.
- N. Nesa, T. Ghosh, and I. Banerjee, "Outlier detection in sensed data using statistical learning models for IoT," in *Proc. IEEE WCNC*, 2018.
- I. Cramer et al., "Detecting anomalies in device event data in the IoT," in *Proc. IoTBDS*, 2018, pp. 52–62