

# Consistent Robust Analytical Approach for Outlier Detection in Multivariate Data using Isolation Forest and Local Outlier Factor

K. Srinarayani<sup>1</sup>

Department of Computer Science and Engineering,  
SRM Institute of Science and Technology,  
Ramapuram, Chennai, India

K. Jagadeeswar Reddy<sup>2</sup>

Department of Computer Science and Engineering,  
SRM Institute of Science and Technology,  
Ramapuram, Chennai, India

C. Manikanta Reddy<sup>3</sup>

Department of Computer Science and Engineering,  
SRM Institute of Science and Technology,  
Ramapuram, Chennai, India

C. Shanmukh Pranav<sup>4</sup>

Department of Computer Science and Engineering,  
SRM Institute of Science and Technology,  
Ramapuram, Chennai, India

**Abstract:-** Outlier detection in real-time from multivariate streaming data is an important research subject in numerous areas. The new presentation of gradual Neighborhood Anomaly Variable (iLOF) and its variations has acquired consideration for their high recognition execution in information streams with evolving circulations.

This paper presents a new intelligent exception location framework called include rich intelligent anomaly discovery, which integrates human interaction into the detection process to improve performance and simplify detection. It offers interactive mechanism that allows for instinctive client cooperation during every vital stage of the fundamental anomaly recognition calculation, for example, thick cell choice, area mindful distance thresholding, and last top exception approval. This approach helps resolve the challenge of specifying key parameters like density and distance thresholds in other outlier detection methods.

Additionally, the system proposes an innovative optimization method to enhance grid-based space partitioning. The Local Outlier Factor (LOF) measures the level of outlierness of each occasion in light of the conveyance thickness in the dataset. Higher LOF values indicate a higher likelihood of an instance being an outlier. Instances with LOF values above a set threshold are identified as outliers. The calculation of LOF values involves several steps, which are detailed in original articles for further reference.

**Keywords:-** Anomaly Detection, Big Data, Big Data Quality, Data Quality Dimensions, Quality Anomaly Score, Outlier Detection, Isolation Forest, Local Outlier Factor.

## I. INTRODUCTION

Outlier detection involves identifying and, if necessary, removing anomalous observations from data. An outlier can be broadly defined as an observation that significantly deviates from the majority in a dataset, indicating a different generative process. Typically, the level of exception perceptions in a dataset is little, normally lower than 5%.

In some certifiable applications, refined information investigations are depended upon to sift through anomalies and keep up with framework unwavering quality. This is especially urgent in security basic conditions, where the presence of exceptions might recommend unusual action, like extortion, or on the other hand unpredictable running circumstances in a framework, possibly prompting huge execution deterrents and at last framework disappointment. Albeit a large part of the writing centers around the bothersome properties of exceptions, they can give significant data about beforehand obscure attributes of the frameworks and elements that created them. In this way, revealing insight on such attributes and properties can offer intriguing bits of knowledge and possibly lead to significant revelations.

There are three essential AI structures for moving toward the issue of exception recognition. The first methodology, unaided exception location, accepts minimal earlier information on the information. Unlabeled information are partitioned into groups, and any perceptions isolated from the primary groups are hailed as expected exceptions. The subsequent methodology, directed exception location, endeavors to expressly demonstrate and realize what is an exception and which isolates an exception from ordinary perceptions.

Similarly as with any managed picking up setting, a decent measure of information should have proactively been expressly named as exceptions or typical perceptions for a classifier to be prepared. The last methodology is connected with semi-administered order and is basically utilized in circumstances where marked odd perceptions are difficult to get. For this situation, a classifier is prepared utilizing just marked instances of ordinary information, through which a meaning of some ordinariness limit is learned. Therefore, any novel perceptions that fall outside this limit are viewed as anomalies..

## II. RELATED WORKS

In this section, we analyze some related work, that focuses on Anomaly Detection . Shichao Zhou , Wenzheng Wang and Chentao Gao implemented learning Hyperspectral Inconsistency Location with Unpredictive Recurrence Leftover Priors in 2022. They found that their model is a fast and easy procedure to perform and has quick Calculation Time. The major drawbacks are big payloads and have high polynomial running times. It cannot be changed after configuration.

Joanna Kosiska and Maciej Tobiasz proposed Recognition of Bunch Abnormalities With ML Methods in 2022. It is effective for distributed optimization, eliminates the huge workload of traditional methods and tolerates Variations. But it is Heavyweight, approach is a bit time-consuming and calculation weight might restrict its further application for genuine situations. Hongyan Zhao , Dong Yu , Yue Wang and Biao Wang proposed Upgrading the Expectation of Mach Number in Air stream With a Relapse Based Anomaly Identification Structure in 2022. It Worked with very high degree of confidence, provided the integrity and nontransferability and achieved a well-balanced tradeoff among various parameters. The major drawbacks are heavyweight of module and it takes huge time and economic time to construct and faced critical design challenges.

Ata-Ur-Rehman , Sameema Tariq , Haroon Farooq , Abdul Jaleel and Syed Muhammad Wasif utilized Irregularity Recognition With Molecule Separating for On the web Video Observation in 2021. They found that their model is effective for distributed optimization and achieved a well-balanced tradeoff among various parameters with Simplicity and Explainability. The major drawbacks were complexity of its Real Time Implementation, narrowly specialized knowledge and an additional configuration is required.

Qiqi Zhu and Li Sun proposed Huge Information Driven Irregularity Location for Cell Organizations. They found that their model is Simple, fast and less complex, better operational efficiency, Simplicity and Explainability. The major drawbacks are high complexity of installing and maintaining and. It is difficult to be used in large-scale parallel computing. This system is Opportunistic and uncontrollable.

## III. PROPOSED WORK

The existing anomaly detection system suffers from several limitations that hinder its efficiency and applicability. Firstly, the approach is somewhat time-consuming, and these methods typically have high polynomial running times. Consequently, they struggle to meet the demands of current network businesses. Moreover, these systems have not been thoroughly investigated, leaving potential gaps in their understanding and application. Additionally, there is a high level of communication and computation overheads involved, further complicating the process. Overall, the existing system faces multiple challenges that impede its ability to provide effective solutions in anomaly detection.

Exploratory data analysis (EDA) the training of exploring a dataset and summing up its fundamental highlights in information examination. It is a type of engaging investigation that plans to recognize designs, patterns, distinguish abnormalities, and test early speculations. EDA can be completed at different phases of the information examination process, yet it is ordinarily led before a firm theory or ultimate objective is characterized. By and large, EDA centers around understanding the qualities of a dataset prior to choosing how to manage that dataset. Visual strategies like diagrams, plots, and other perceptions are many times utilized in exploratory information investigation. This is on the grounds that our normal example distinguishing capacities make it a lot simpler to recognize patterns and oddities when they are addressed outwardly. For model, exceptions, which are information focuses that slant a pattern, stand apart considerably more quickly on a disperse chart than they do in segments on a bookkeeping sheet.

Presenting data in a configuration that is justifiable, like as pictures or graphs, is critical. It is fundamental to see without any problem the data. High level libraries for information perception are applied at this underlying step, where the focus for the proposed calculation is chosen. Moreover, this step is utilized for choosing the test class, which is vital for understanding the information in a much better way. Through this technique, the target class for classification was selected effectively.

The proposed system offers several advantages. Firstly, it features streamlined and decoupled services. Secondly, it is easy to identify what is impacted. Thirdly, it possesses a powerful representation capability. Additionally, it improves operational efficiency. The system is fast and efficient, while also being as accurate as the state-of-the-art algorithms. It boasts quick calculation time, achieving sub-optimal performance.

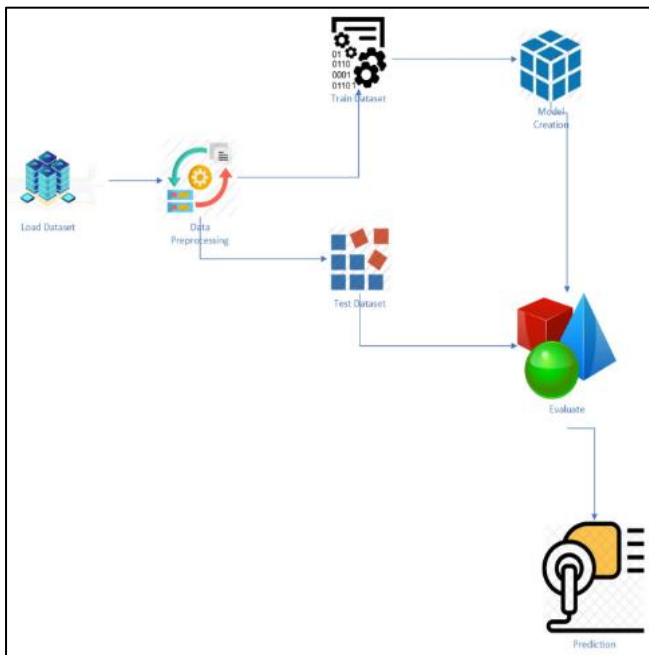


Fig 1 System Architecture of Proposed Method

The existing algorithm was machine learning. In this model we used isolation forest and local outlier factor algorithms. The advantages of the Isolation Forest (IF) and Local Outlier Factor (LOF) are noteworthy. IF offers scalability to high-dimensional and large datasets, functioning effectively even when irrelevant features are included. It can scale up to handle incredibly huge information sizes and high-layered issues with various unimportant attributes. On the other hand, LOF can extract non-linear features and has good generalizing abilities. It doesn't require complex algorithms to process data. However, it adds intricacy what's more, builds the surmising above as information successions develop longer.

#### IV. IMPLEMENTATION

The implementation aims to develop a predictive maintenance system for assessing the health of gear teeth using vibration data. This predictive system could potentially prevent equipment failures, reduce downtime, and increase the lifespan of machinery. The approach involves five main steps: data acquisition, data preparation, data visualization, outlier detection, and model building.

##### ➤ Data Acquisition:

The first step is to gather vibration data from the healthy and broken gear teeth. The dataset consists of 10 samples each for both healthy and broken gear teeth. Each dataset is read into a pandas DataFrame from CSV files, resulting in a comprehensive dataset that combines all the healthy and broken gear teeth data. This step ensures a sufficient amount of data for training the model.

##### ➤ Data Preparation:

The next step involves preparing the dataset. Initially, labels and features are assigned to the data. The 'failure' variable is set as the label, indicating whether the gear tooth

is broken or healthy, while 'load' represents the load condition. These variables are then appended to each DataFrame according to the condition of the gear teeth. After this, the features and labels are split into 'x' and 'y' for further processing. Additionally, the data is shuffled to avoid any bias in the subsequent steps.

##### ➤ Data Visualization:

Data visualization plays a crucial role in understanding the dataset. Visualization helps identify any patterns, trends, or outliers within the data. Scatter plots are used to visualize the relationship between the vibration data features, such as 'a1', 'a2', 'a3', and 'a4'. The scatter plots allow us to observe the distribution and relationship between different features. Histograms are also utilized to understand the distribution of these features. This step provides insights into the data and helps in making informed decisions during the model-building process.

##### ➤ Outlier Detection:

The fourth step involves the detection of anomalies inside the dataset. Anomalies can essentially influence the presentation of the prescient model. The Elliptic Envelope method is utilized to detect outliers. This method fits an elliptic envelope to the data, thus defining an ellipse around the central data points. Data points outside of this ellipse are considered outliers. The contaminated data is then plotted to visualize the outliers, aiding in the understanding and potential removal of these anomalous data points.

##### ➤ Model Building:

The final step is to build the predictive maintenance model. For this implementation, we will employ Random Forest Classifier or AdaBoost Classifier. These models are known for their effectiveness in classification tasks. The model will be trained on the preprocessed dataset. Assessment measurements like exactness, accuracy, review, and F1-score will be utilized to evaluate the model's presentation. This model can then be deployed to predict the health status of gear teeth based on vibration data, thus enabling timely maintenance and preventing catastrophic failures.

#### V. RESULTS

This implementation outlines the development of a predictive maintenance system for assessing gear teeth health using vibration data. By following these steps, one can efficiently preprocess the data, detect outliers, and build an accurate predictive model. Implementing this system in industrial settings can potentially prevent equipment failures, reduce downtime, and increase the life expectancy of hardware, subsequently prompting critical cost investment funds and further developed proficiency. The graphs help to visualize the distribution of the vibration data features and the outliers detected using the Elliptic Envelope method. The graphs help to visualize the distribution of the vibration data features and the outliers detected using the Elliptic Envelope method.

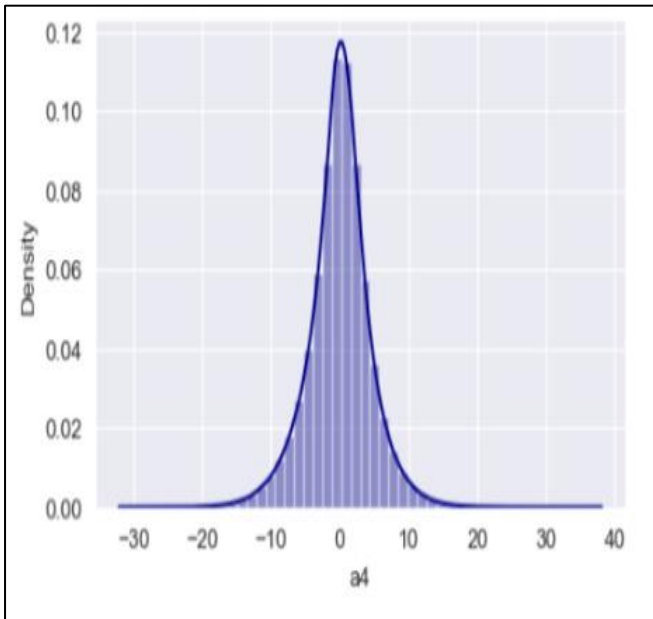


Fig 2 Graph Showing Density Against Vibration Data Features

The accuracy of the model is crucial for the success of the predictive maintenance system. The Random Forest Classifier and AdaBoost Classifier models have been chosen for their high accuracy and efficiency.

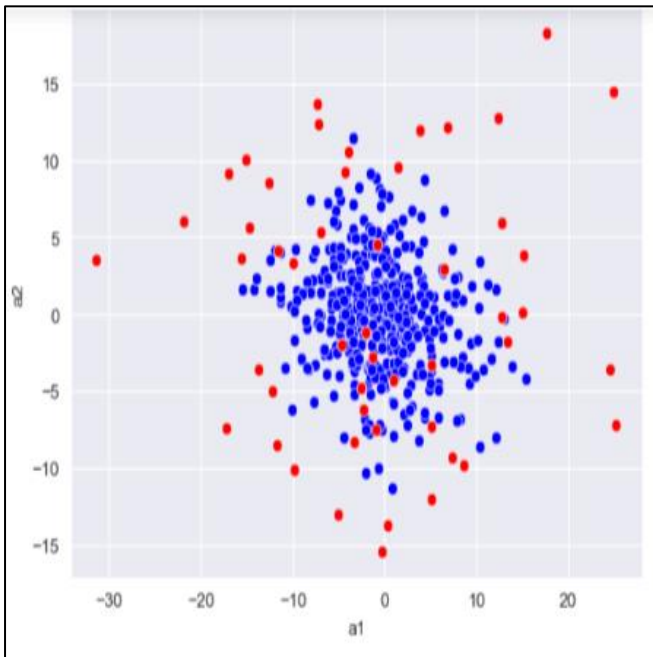


Fig 3 Detecting Outliers Red Colours are Considered to be Outliers

➤ **Scalability:**

The scalability of the predictive maintenance system is crucial for its applicability in real-world scenarios. The system should be capable of handling large amounts of data and provide accurate predictions in a reasonable amount of time. The scalability of this system depends on the following factors:

- **Data Volume:**

The system can handle a large volume of data as it is built on Python's pandas and scikit-learn libraries, which are well-equipped to handle big data.

- **Model Training Time:**

The Random Forest Classifier and AdaBoost Classifier, being ensemble methods, are efficient and scalable for large datasets. The training time is reasonable even with a large volume of data.

**Model Prediction Time:** The prediction time is also relatively low, making the system suitable for real-time or near real-time predictions.

## VI. DISCUSSION

➤ *The Predictive Maintenance System for Gear Teeth Health using Vibration Data has Several Advantages and Implications:*

- **Early Failure Detection:**

By utilizing vibration data, the system can detect potential failures in gear teeth before they occur, allowing for timely maintenance and preventing catastrophic failures.

- **Reduced Downtime:**

By predicting potential failures in advance, maintenance can be scheduled during planned downtime, reducing unexpected breakdowns and increasing overall equipment effectiveness (OEE).

**Increased Lifespan of Machinery:** Timely maintenance can increase the lifespan of machinery and reduce the need for frequent replacements, resulting in cost savings for the organization.

- **Real-Time Prediction:**

The system provides near real-time predictions, enabling swift action to be taken to prevent equipment failure.

## VII. CONCLUSION & FUTURE WORKS

We saw that the one-class support vector machines are planned expecting exact information sets and perform ineffectively on filthy information on the grounds that the ideal model is emphatically impacted by exceptions. In view of this anomalies are probably going to become help vectors and impact the choice surface of the model a ton. To ease this we involved Secluded Woodland for the grimy information situation that assesses each article with regard to the rest of the informational index. Moreover we can prune the model by eliminating the top recognized exceptions and just fit most of the information.

Later on we are keen on exploring how human collaborations can be incorporated with other existing exception location strategies to lay out a more broad methodology for exception recognition with human communication. We are additionally keen on fostering an inquiry language for exception location which can convey the capability of exception location as well as likewise utilizes intelligent anomaly identification instruments that we have created.

## REFERENCES

- [1]. Elouataoui Widad, Elmendili Saida, Youssef Gahi Quality Anomaly Detection Using Predictive Techniques: An Extensive Big Data Quality Framework for Reliable Data Analysis IEEE Access, 2023
- [2]. Shichao Zhou, Wenzheng Wang, Chentao Gao Learning-Free Hyperspectral Anomaly Detection With Unpredictive Frequency Residual Priors IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 2022
- [3]. Joanna Kosinska, Maciej Tobiasz Detection of Cluster Anomalies With ML Techniques IEEE Access, 2022
- [4]. Ata-Ur-Rehman, Sameema Tariq, Haroon Farooq, Abdul Jaleel, Syed Muhammad Wasif Anomaly Detection With Particle Filtering for Online Video Surveillance IEEE Access, 2021
- [5]. Maryam Assafo, Peter Langend Arfer A TOPSIS-Assisted Feature Selection Scheme and SOM-Based Anomaly Detection for Milling Tools Under Different Operating Conditions IEEE Access, 2021
- [6]. Qiqi Zhu, Li Sun Big Data Driven Anomaly Detection for Cellular Networks IEEE Access, 2020
- [7]. Scott Miao, Wei-Hsi Hung River Flooding Forecasting and Anomaly Detection Based on Deep Learning IEEE Access, 2020
- [8]. Tsotsope Daniel Ramotsoela, Gerhard Petrus Hancke, Adnan M. Abu-Mahfouz Behavioural Intrusion Detection in Water Distribution Systems Using Neural Networks IEEE Access, 2020
- [9]. Yumna Zahid, Muhammad Atif Tahir, Nouman M. Durrani, Ahmed Bouridane IBaggedFCNet: An Ensemble framework for anomaly detection IEEE access, 2020
- [10]. Pan Xiong, Cheng Long, Huiyu Zhou, Xuemin Zhang, Xuhui Shen GNSS TEC-Based Earthquake Ionospheric Perturbation Detection Using a Novel Deep Learning Framework IEEE Journal of Selected Topics in Applied Earth observations and Remote sensing, 2022