# Uncovering the Hidden Patterns: A Novel Approach to Anomaly Detection with Machine Learning

**Deepti Negi**

Asst. Professor, School of Computing, Graphic Era Hill University, Dehradun, Uttarakhand
India 248002

## ABSTRACT

Anomaly detection has a distinct level of issue complexity than the vast majority of analytical and learning issues and activities. This is due to the singular nature of the phenomenon being investigated. The aforementioned inherent obstacles and unresolved identification concerns in complex anomaly data are compiled in this paper. It has come to the attention of those working in data mining, machine learning, computer vision and statistics that the process of anomaly detection is becoming an increasingly important one. This is a result of the increasing demand for the application of these technologies across a wide range of sectors, including risk management, security, compliance, financial monitoring, and AI safety. Deep learning has lately pushed the bounds of various different learning tasks due to its incredible ability to learn expressive representations of complicated data, such as high-dimensional data, temporal data, geographical data, and graph data. Deep learning for anomaly detection, also known as deep ADE, is an application of deep learning that makes use of neural networks to learn feature representations or anomaly scores in order to identify anomalies. When it comes to resolving difficult detection problems in a range of real-world circumstances, several deep anomaly detection algorithms have been developed, and they have demonstrated to perform noticeably better than standard anomaly detection techniques. This article will attempt to provide an all-encompassing analysis of the subject matter covered here.

## 1. INTRODUCTION

The anomaly detection approach is used to locate a strange point or pattern in a collection of data. Sometimes the term "anomaly" is also referred to as "outlier." The data mining specialists had previously focused on other techniques like categorization and clustering. Outliers are found during the data cleansing process. The impression was altered in 2000, however, when researchers found that the identification of aberrant goods may help with problems with damage detection, fraud detection, intrusion detection, and the detection of abnormal health conditions. Contextual anomalies, point anomalies and collective anomalies are the three types of anomalies that are included under the phrase "anomalies". When one instance in a data-set deviates from the rest in terms of one or more of its attributes, this is known as a point anomaly.

Anomaly detection (also known as outlier detection) is the process of identifying unusual items, occurrences, or observations that raise concerns by departing dramatically from the rest of the data.

Atypical data is often linked to a problem or unexpected event, such as, for instance, financial fraud, medical problems, structural problems, malfunctioning equipment, etc. Knowing which data points may be classified as anomalies due to this link is especially intriguing since recognizing these occurrences is often quite important from a business viewpoint.

### *Anomaly Detection Approach*

Methods for detecting anomalies can be classified according to the kinds of data that are required to train the corresponding models. It is reasonable to anticipate that outlier samples will only constitute a very small part of the overall dataset in the vast majority of practical applications. As a result, normal data samples are far easier to identify than aberrant data samples, and this is true even when labelled data is available. This assumption is necessary for virtually all applications in the modern world. In the following sections, we will discuss how the availability of labelled data influences the approach that is selected and how this will be done.

1) Supervised Learning:

Machines can learn a function that translates input characteristics into outputs when they are taught using examples of input-output pairs when they are learning in a supervised setting. The incorporation of information that is unique to the application being used into the process of anomaly identification is the objective of supervised anomaly detection algorithms. It is possible to recast the challenge of anomaly detection as a classification task by providing sufficient instances of normal and abnormal behavior. This will allow computers to learn how to accurately anticipate whether or not a given example demonstrates aberrant behavior. Although there may be multiple kinds of anomalies, each one may be greatly underrepresented in many anomalies' detection use cases where the ratio of normal to abnormal occurrences is drastically skewed.

This approach assumes that the user has examples of each kind of anomaly and can appropriately categorize all conceivable abnormalities. This is often not the case in practice since abnormalities may present in a number of ways and additional anomalies can develop during testing. In order to find anomalies that haven't been noticed previously, approaches that generalize effectively are recommended.
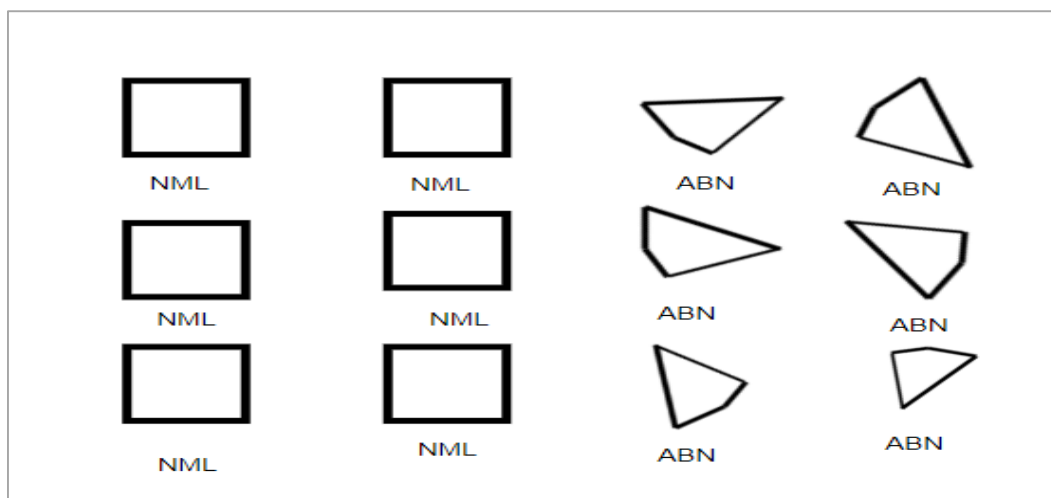


Figure 1. Supervised Learning Training

2) Unsupervised Learning:

Because there aren't enough instances of input-output pairs, machines cannot learn a function that converts input characteristics to outputs via unsupervised learning. They instead find structure in the input information and make use of it to learn. Since tagged anomalous data is, as was previously said, rather infrequent, unsupervised approaches are more often utilized than supervised ones in the area of anomaly identification. The abnormalities that one expects to detect, however, are typically highly specific. Because of this, many irregularities found using a completely unsupervised technique could just be noise and have no bearing on the work at hand.

3) Semi Supervised Learning:

As a kind of middle ground, semi-supervised learning algorithms utilize a number of approaches that can profit from both massive volumes of unlabelled data and sparsely labelled data. Many real-world anomaly detections use cases are well-suited to semi-supervised learning because there are so many examples of the more uncommon or abnormal classes of interest and so few examples of the more normal classes from which to learn. A robust model may be trained on an unlabelled dataset with the assumption that the majority of the data points are normal, and its performance can be evaluated (along with the model's parameters) using a modest amount of labelled data.

This hybrid technique is most suited for applications such as network intrusion detection where there may be several instances of the normal class and a few examples of intrusion classes, but new types of intrusions may arise over time. In other words, there may be numerous instances of the normal class and a few examples of intrusion classes. One such illustration would be the use of X-ray screening for safety purposes at airports and borders. Products that could compromise one's safety are extremely uncommon and can take on a number of different shapes. A potential hazard's anomaly may also vary in nature as a consequence of a range of external occurrences. Consequently, it could be difficult to gather enough meaningful instances of anomaly data in large enough amounts.

## 2. LITERATURE SURVEY

Chahla, Charbel et al [1] studied anomalies in time series data has great promise for reducing energy waste and improving building energy consumption monitoring. For the purpose of identifying outliers in time series data, they described a hybrid model in this study that combines the LSTM and K-means algorithms. While the suggested approach locates the average scenario and may localise the identified anomalies, Autoencoders detect atypical days. Despite these positive findings, their effort still needs the input of genuine professional users and analysts in order to more precisely identify the abnormality in this field. Experts can also add annotations to the learning data, giving them the opportunity to use semi-supervised techniques in this field.

Chandola Varun et al [2] in their review, they have discussed numerous formulations of the topic of anomaly detection in the literature and have made an effort to give a general overview of the vast literature on various methodologies. they have selected a specific premise about what constitutes normal and abnormal data for each category of anomaly detection tools. These suppositions can be used as guides to calculate whether a actual technique is effective in a given area when applied to it. A thorough examination of anomaly detection should ideally enable readers to not only comprehend the rationale for selecting a specific technique but also to compare and contrast alternative strategies. However, the current research has been conducted in an unorganised manner and

without relying on a common concept of anomalies, which makes it exceedingly difficult to provide a theoretical knowledge of the anomaly detection problem.

A novel paradigm for successful early detection and identification of network abnormalities across high-speed networks is presented in the paper by Osman Salem et al. [3], enabling quick response and the execution of countermeasures. The recommended architecture uses change point detection in the counter values of reversible drawings to gather multiple data streams from high-speed lines in a stretched database. They use a unique technique for sketch inversion to discover the occurrence of change points and show the culprit flows in order to find network anomalies. They use the cumulative sum (CUSUM) algorithm at the counter value of each bucket in the recommended reversible drawing. A theoretical framework for attack detection is presented. As part of the OSCAR French research project, they also provide the results of their studies, which were carried out on two actual data trails with abnormalities and carefully evaluated.

Pang Guansong et al [4] evaluated 12 different modeling perspectives on the use of deep learning techniques for anomaly detection. They also spoke about how these strategies address a number of well-known anomaly detection challenges in order to emphasize the importance of deep anomaly detection.

Their book covers current research on representative approaches to demonstrate how a data-centric system designed by Arun Kumar Sangaiah et al. [5] may be used to analyze, process, and describe cloud-hosted multimedia massive data utilizing CI. Modeling, relationship pattern identification, clustering, and other bioengineering challenges may be solved using CI methods, according to the book. Programmers and subject matter experts who want to understand and research the use of computational intelligence aspects (opportunities and challenges) for the

design and development of a data-centric system in the multimedia cloud, big data era, and its related applications, such as smarter healthcare, traffic control, homeland security, trading analysis, and telecom, are the target audience. Researchers and PhD candidates studying data-centric systems' role in the future of computing can benefit from this book.

A novel LSTM-based technique (LSTMAD) for detecting abnormalities in time series data was proposed by Ji Zhiwei et al. [6]. LSTMAD was developed by merging the LSTM network with a statistical technique. Since their technique uses prediction error to detect abnormal regions after learning the context of sequence data from the normal signal, there is no need to depend on previous information. To test the performance, they used LSTMAD to a variety of time series datasets, including well-known public data and real-world data. The results of the simulation demonstrated that LSTMAD is capable of accurately identifying differences in an entire sequence. LSTMAD performed better than the other gold standard methods on every testing dataset.

Ahn Jaesung et al [7] provide a unique approach called Multi-Level Masking and Restoration with Refinement 54 (MMRR), which overcomes the hyperpa55 remoter sensitivity issue without the use of side information and is based on a generative model. Their suggested strategy is motivated by the idea that a network that has been trained to recover normal data from scant information about normal data would also learn the prominent elements of that data. Anomaly identification in terms of restoration can be done since restoration from limited information succeeds for normal data and fails for 58 aberrant data. The following two elements are the main parts of their method in order to achieve this goal. First, masking, which is a technique 60 that restricts the remaining information with the exception of 61 the sections

required for restoration in order to get restricted information. The second step is restoration, which is the act of regaining access to the original data 62 utilising just the constrained data gleaned by masking. 63 Finding the ideal masking level—64, which is the degree to which the mask limits information—is important for MMRR to conduct anomaly detection. At this level, normal data can be successfully restored but abnormal data cannot.

However, rather than locating a single optimal masking level, they employed 67 ensembles at various masking levels to detect anomalies in order to circumvent the hyperparameter sensitivity 66 issue brought on by the lack of aberrant data during training. By allowing the 69 manual control of the mask's masking level, their new 68 mask generation method enabled ensemble at different masking levels, eliminating the necessity for adversarial loss. Additionally, their approach of mask synthesis made the mask learnable such that the mask that was most useful for restoration at the corresponding masking level was formed, which improved the performance of anomaly detection. 73 their masking method, however, compares the restoration level at the same masking level without taking the complexity of each data set into account. As a result, basic anomalous data can typically be restored using masking and restoration alone rather than complex normal data, in which case anomaly detection 76 often unsuccessful. To address this issue, they suggested an additional refining method that eliminates the variation in restoration brought on by the variation in data complexity.

An innovative adversarial training architecture for unsupervised anomaly recognition is presented in S. Akçay et al.'s [8] study. The recommended technique looks at how skip connections work within the generator and feature extraction from the discriminator to alter hidden features. The findings demonstrate that skip connections provide more stable training

and that inference learning from the discriminator outperforms earlier state-of-the-art techniques quantitatively. This is based on an analysis of a large number of datasets with varying levels of complexity and domain. The study's actual findings provide insight into how well the recommended technique may be applied to any task involving the discovery of anomalies. It may be worthwhile to look into the recommended method's applicability to higher quality images and other temporal information-based anomaly identification jobs.

Bergman Liron et al [9] for three datasets, they ran contamination experiments. Due to the insufficient quantity of anomalies, thyroid was left out. In contrast to KDDRev, their methodology does not employ unused anomalies for contamination. The anomalies are instead divided into train and test groups. Anomalies in tests are utilised for evaluation, while anomalies in trains are used for contamination. They only present GOAD because DAGMM did not present findings for the other datasets. In terms of contamination on KDD, KDDRev, and Arrhythmia, GOAD was comparatively resilient. They provided a technique for identifying irregularities in general data. This was accomplished by using a collection of auxiliary tasks to train a classifier. They are able to create an infinite number of random tasks using their proposed method, that not really need any former knowledge of the data domain. Their approach vastly outperforms the state-of-the-art.

T. Chen et al. [10] presented a basic contrastive learning paradigm for visual representations. They simplify new contrastive self-supervised learning algorithms without specific structures or a memory bank. They examine the framework's core components to determine what allows useful representations in contrastive prediction challenges. They demonstrate data augmentations for

efficient prediction tasks, the addition of a learnable nonlinear transformation between the learned representation and the contrastive loss, and the advantages of contrastive learning over supervised learning in terms of larger batch sizes and more training steps. They significantly outperformed self-supervised and semi-supervised learning approaches on ImageNet by integrating these findings. A linear classifier trained on SimCLR-learned self-supervised representations matches a supervised ResNet-50 and achieves 76.5% top-1 accuracy, 7% better than the previous state-of-the-art. Their model achieves 85.8% top-5 accuracy with 1% fine-tuning, surpassing AlexNet with 100 fewer labels.

## 3. PROPOSED SYSTEM

Deep learning algorithms offer several advantages, which is why we utilized them to identify anomalies. These approaches are designed to work with multivariate and highly dimensional data, first and foremost. This makes it straightforward to integrate data from several sources since it eliminates the challenges of independently modelling anomalies for each variable and averaging the results. Deep learning algorithms are particularly adapted to simultaneously modelling the interactions between numerous variables with regard to a specific task, in addition to setting general hyperparameters (number of layers, units per layer, etc.). Deep learning models only need modest adjustments to provide worthwhile outcomes.

Another advantage is performance. The ability to model complex, nonlinear connections within data and utilise these interactions for the job of anomaly detection is made possible by deep learning methods. Because their effectiveness may scale with the quantity of relevant training data, deep learning models are appropriate for problems involving large amounts of data.
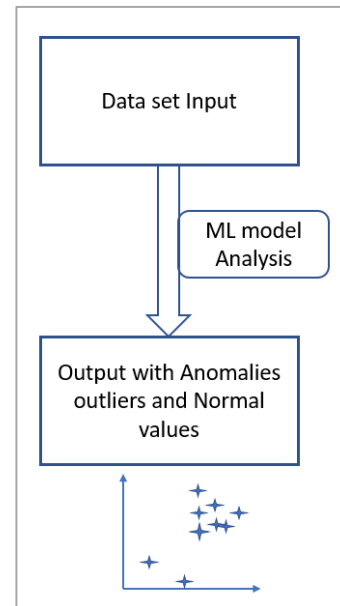


Figure 2. Proposed ML model for Anomaly detection

### A. Local Outlier Factor

Local outlier factor is perhaps the most often used technique for anomaly detection. This approach is based on the notion of the local density. It compares the local density of an item to the densities of the close-by data points. If a data point's density is lower than that of its neighbors, it is considered an outlier.
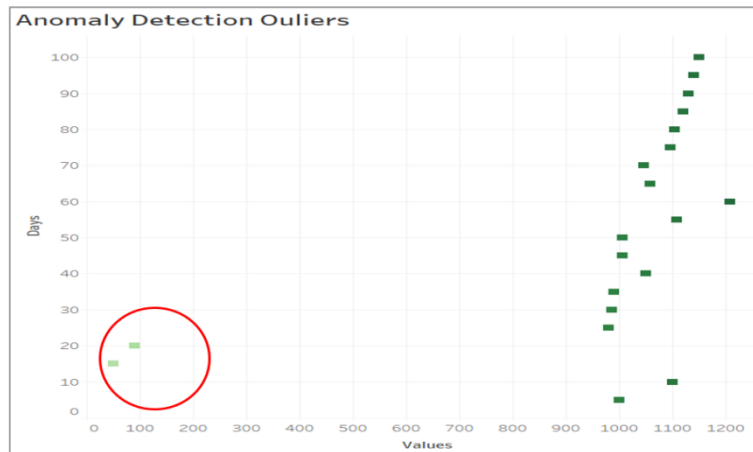
Figure 3. Anomaly Detection with scatter Plot analysis

In above figure we can see we have used dataset which contains values proportional to days. Where marked with red circle states Q1 that those are anomalies detected as outliers and other set nearest data values Q2 are grouped together as they shared nearby properties n values.

We can see pattern where suddenly few data points get very odd or we can see outliers' values where as others are still intact to near by points.

Formulas for outliers Lower and Higher Outlier as below:

Lower Outlier values (Lower Bound) = Q1 – (1.5 * IQR)
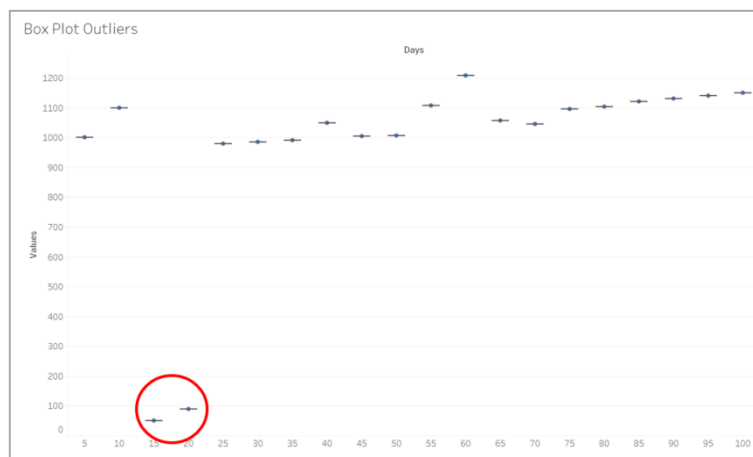Higher Outlier values (Higher Bound)= Q3 + (1.5 * IQR)



Figure 4. Box Plot analysis for Outliers Anomaly Detection

## 4.  CONCLUSION

We established set up where we used dataset to check outliers using LOF techniques with supervised method. We used Q1 and Q2 set to differs between normal (Higher outliers) and abnormal data (lower outliers) values and detected lower Q1 set as an outlier whereas Q2 as higher bound of outliers which have similar types of values. We also analysed data using scatter plot chart and box plot chart which both gave results which confirmed Q1 lower bound have some distinct values where Q2 higher bound have another set of values.

**REFERENCE**

[1] Chahla, Charbel, Hichem Snoussi, Leïla Merghem, and Moez Esseghir. "A Novel Approach for Anomaly Detection in Power Consumption Data." In *ICPRAM*, pp. 483-490. 2019.

[2] Chandola, Varun, Arindam Banerjee, and Vipin Kumar. "Anomaly detection: A survey." *ACM computing surveys (CSUR)* 41, no. 3 (2009): 1-58.

[3] Salem, Osman, Sandrine Vaton, and Annie Gravey. "A Novel Approach for Anomaly Detection over." *on Computer Network Defense* (2010): 49.

[4] Pang, Guansong, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. "Deep learning for anomaly detection: A review." *ACM computing surveys (CSUR)* 54, no. 2 (2021): 1-38.

[5] Books: Computational Intelligence for Multimedia Big Data on the Cloud with Engineering Applications. Editors: Arun Kumar Sangaiah, Zhiyong Z,Zhang, Michael Sheng. Paperback ISBN: 9780128133149, 2018

[6] Ji, Zhiwei, Jiaheng Gong, and Jiarui Feng. "A novel deep learning approach for anomaly detection of time series data." *Scientific Programming* 2021 (2021).

[7] Ahn, Jaesung, Janghyeon Lee, Hanbyel Cho, Yooshin Cho, HyeongGwon Hong, and Junmo Kim. "MMRR: Unsupervised Anomaly Detection through Multi-Level Masking and Restoration with Refinement.", 2021

[8] S. Akçay, A. Atapour-Abarghouei, and T. P. Breckon. Skip-ganomaly: Skip connected and adversarially 379 trained encoder-decoder anomaly detection. In 2019 International Joint Conference on Neural Networks 380 (IJCNN), pages 1–8. IEEE, 2019.

[9] Bergman, Liron, and Yedid Hoshen. "Classification-based anomaly detection for general data." *arXiv preprint arXiv:2005.02359* (2020).

[10] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual 394 representations. In International conference on machine learning, pages 1597–1607. PMLR, 2020