

Adaptive Thresholding Heuristic for KPI Anomaly Detection

Ebenezer R.H.P. Isaac and Akshat Sharma
Global AI Accelerator, Ericsson, Chennai, India
ebeisaac@ieee.org, akshat.vhs@gmail.com

Abstract—A plethora of outlier detectors have been explored in the time series domain, however, in a business sense, not all outliers are anomalies of interest. Existing anomaly detection solutions are confined to certain outlier detectors limiting their applicability to broader anomaly detection use cases. Network KPIs (Key Performance Indicators) tend to exhibit stochastic behaviour producing statistical outliers, most of which do not adversely affect business operations. Thus, a heuristic is required to capture the business definition of an anomaly for time series KPI. This article proposes an Adaptive Thresholding Heuristic (ATH) to dynamically adjust the detection threshold based on the local properties of the data distribution and adapt to changes in time series patterns. The heuristic derives the threshold based on the expected periodicity and the observed proportion of anomalies minimizing false positives and addressing concept drift. ATH can be used in conjunction with any underlying seasonality decomposition method and an outlier detector that yields an outlier score. This method has been tested on EON1-Cell-U, a labeled KPI anomaly dataset produced by Ericsson, to validate our hypothesis. Experimental results show that ATH is computationally efficient making it scalable for near real time anomaly detection and flexible with multiple forecasters and outlier detectors.

Index Terms—Pattern recognition, statistical learning, time series, unsupervised learning, telecom AI.

I. INTRODUCTION

Instances that significantly deviate from the observed statistical pattern in the data are called outliers. Outlier detection is a crucial step in any data science application. Generally, the terms outliers and anomalies are used interchangeably. Nevertheless, within the realm of business, an anomaly can be defined as an unexpected event or situation that pertains to a particular use case. All anomalies may be considered outliers, but not all outliers are anomalies. Network KPI anomaly detection (AD) is an integral part of many use cases in the telecom domain including sleeping cell detection [1], ensuring SLA adherence [2], and abnormal traffic detection [3]. AD has become increasingly important as telecom data continue to grow at an exponential rate. Outlier detectors are usually compared in terms of the area under the receiver operating characteristic curve (AUC-ROC), the greater the area, the better is the tradeoff between true and false positives. However, it is necessary to understand how to set the threshold for a given use case and that this threshold may not be fixed for the life-cycle of the use case. At times, the traffic pattern can change altering the definition of the norm and hence requiring recomputation of the threshold.

Recently, AD solutions based on Deep Learning (DL) have been gaining popularity [4], [5], [6], [7]. However, DL methods require extensive computation load on the system. In a typical telecom use case involving cell KPI datasets, the volume of data processed can range from gigabytes to terabytes per month depending on the number of cells which ranges from thousands to tens of thousands. In such cases, employing a DL-based solution would not be viable. It would hence be better to include standard statistical learning or machine learning methods, many of which are described in the PyOD toolbox [8], and then proceed with a use-case-specific filter to mitigate false positives.

Thresholding in AD defines a limit or boundary to distinguish normal behavior from anomalous behavior in a dataset. A typical $K\sigma$ deviation method involves setting a threshold that is K times the observed standard deviation, σ , from the mean [9] based on the assumption that the data is normally distributed. However, this assumption does not hold true for most telecom data. The annual maximum method identifies anomalies based on extreme values observed within a specific time period, usually a year [10]. It is particularly useful for capturing rare events across longer periods, but not be suitable for detecting anomalies occurring at smaller time scales. LSTM-NDT [11] is a deep learning-based method that employs Long-Short-Term-Memory (LSTM) neural networks with non-parametric dynamic thresholding. While LSTM-NDT can effectively capture complex temporal patterns, it may not scale well to large datasets as it is computationally expensive. The Peak Over Threshold (POT) [12] method involves fitting the data to a Generalized Pareto Distribution to determine threshold values. Tuli et al. [13] employed POT and observed better results in comparison to AM. However, POT itself requires a threshold as input to provide a threshold for anomaly detection which appears to be its main disadvantage.[10].

The most successful thresholding methods in literature are computationally intensive, highly sensitive to non-trivial hyperparameter tuning, and few address concept drift. These limitations makes it difficult for existing methods to balance between false positives and false negatives. The method proposed in this article aims to address the above gaps by selecting a threshold that rules out periodic patterns of statistical outliers with an additional mechanism to address concept drift making it robust against noise. The contribution of this article are summarized as follows.

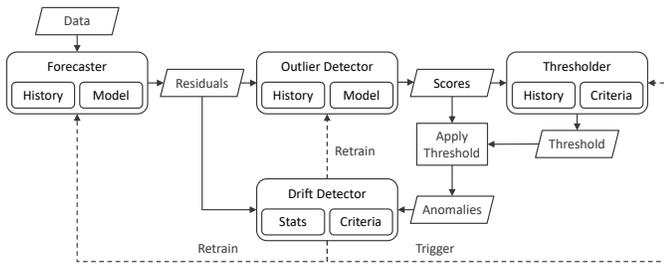


Fig. 1. Simplified flow of time series anomaly detection with ATH

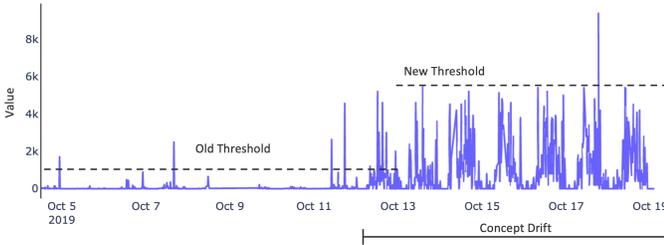


Fig. 2. Concept drift illustration with ATH. The occurrence of a concept drift breaks the ATH criteria triggering a recomputation of the threshold.

- 1) A time series AD thresholding heuristic specialized for near real time KPI anomaly detection interoperable with any forecaster and any outlier detector that returns a score.
- 2) A business logic to differentiate between outliers and use-case-specific anomalies for telecom KPI datasets.
- 3) Introduction of a labelled benchmark dataset for time series AD with telecom KPIs

II. METHOD

This article introduces Adaptive Thresholding Heuristic (ATH), a novel method for setting dynamic thresholds for based on the statistical properties of the data (filed as a patent by Ericsson in [14]). It is particularly well-suited to data with unknown or varying distributions. One can assume that anomalous occurrences in time series data does not follow a periodic pattern and are rare in occurrence. In other words, if there is a set of occurrences that seem to be outliers, and if those occurrences occur periodically or too frequently, then the occurrences are most possibly not anomalous. A heuristic is set to adhere to this assumption as constraints. The threshold of a detector is tuned appropriately such that these constraints are satisfied. The occurrence of a concept drift can be identified should the periodicity or anomaly proportion constraint be broken during operation after the threshold is set. Should a concept drift occur, then, depending on the business use case, a simple threshold update can be made by running this routine again, or, in the worst case (as in multivariate systems), the detectors are retrained, and then the heuristic is reapplied.

An overview of the method is shown in Fig. 1. The initial step for time series AD involves the application of forecasting techniques to obtain the residuals. Residuals represent the deviation between the predicted values and the actual

Algorithm 1 Adaptive Thresholding Heuristic

Input:

X : signal with time series data
 M : threshold-based outlier detector model
 $tail$: either “left” or “right”
 $periodicity_limit$: number of permitted periodic outlier occurrences
 $proportion_limit$: permitted proportion of outliers

Output: threshold

```

1: procedure APPLYATH( $X, M, tail, periodicity\_limit, proportion\_limit$ )
2:   Fit the model  $M$  using  $X$ 
3:    $S \leftarrow$  scores of  $X$  using  $M$ 
4:    $thresh\_list \leftarrow$  unique values of  $S$ 
5:   Sort  $thresh\_list$  by its values according to the specified  $tail$ 
6:     left: ascending order
7:     right: descending order
8:    $previous\_thresh \leftarrow$  first item in  $thresh\_list$ 
9:   for  $thresh$  in  $thresh\_list$  do
10:     $outliers \leftarrow$  empty list
11:    Based on the  $tail$ , do the following
12:     left: add to  $outliers$  all  $s \in S$  such that  $s < thresh$ 
13:     right: add to  $outliers$  all  $s \in S$  such that  $s > thresh$ 
14:     $diff \leftarrow$  list temporal differences between each value pair in  $outliers$ 
15:    Remove the differences within each consecutive outlier from  $diff$ 
16:    if any of following conditions are met then break from the loop
17:     Frequency count of any value of  $diff > periodicity\_limit$ 
18:      $|outliers|/|X| > proportion\_limit$ 
19:    end if
20:     $previous\_thresh \leftarrow thresh$ 
21:  end for
22:   $final\_thresh \leftarrow previous\_thresh$ 
23:  return  $final\_thresh$ 
24: end procedure

```

observations. These residuals are passed through an outlier detector to obtain anomaly scores. Subsequently, thresholding is employed on these scores, whereby the ATH algorithm is utilized. Each stage of the pipeline holds the history of inputs pertaining to a moving window. In a live deployment, the size of the windows may differ from one stage to another. For instance, the forecaster can have an input window of one month while the outlier detector and thresholder can have a window of one week. The drift detector has access to the historical statistics of the percept history of each stage and can trigger a threshold recomputation and retraining should a concept/data drift occur.

A complete description of ATH is provided in Algorithm 1. Given a time series data stream, a detector is fit to a window and the outlier scores are extracted. Then, the scores are sorted and checked until the periodicity condition or the anomaly proportion limit is met. Once, the threshold that breaks the constraints is reached, the previous (most recent) threshold in the list that satisfies the constraints is selected as the final threshold. The selected threshold is used to label the anomalies. Both periodicity and proportionality constraints are also monitored through the drift detector so that ATH can be triggered should any of the conditions fail and hence minimize false positives by increasing the threshold accordingly. Also, the drift detector maintains the residual statistics such as its operating range. For instance, reduction in the operating ranges should also trigger ATH to bring down the threshold reduce false negatives. ATH can be applied to any outlier detector that gives a outlier score as output including a simple Z-Score model.

To check if there is periodicity in the outliers, the system

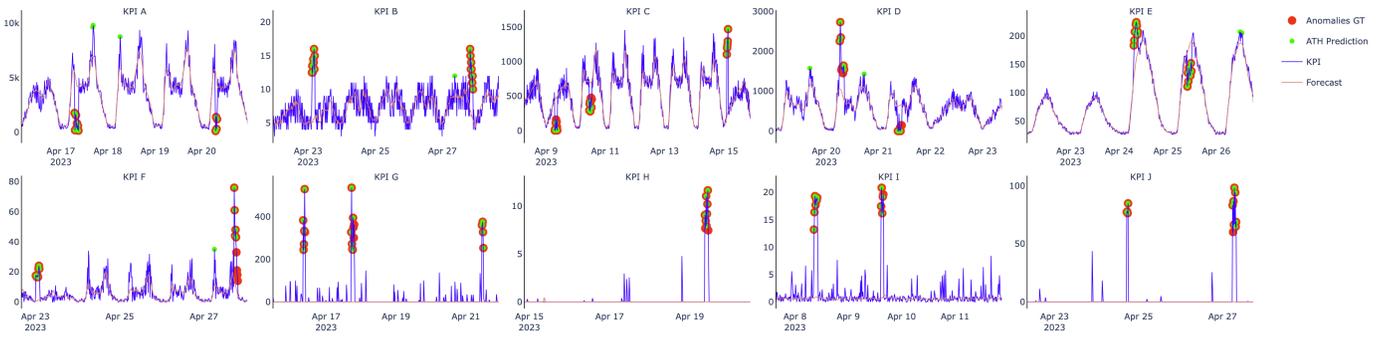


Fig. 3. Illustration of ATH-based AD through QBSD forecaster and Z-Score detector on the EON1-Cell-U dataset. Only subsets are shown for readability.

computes the temporal difference between each outlier. This operation can be done in several ways, depending on the use case. For the KPI data, day-wise differences can be considered. For example, if an outlier occurs every day (regardless of the time of day) more than the *periodicity_limit*, then the observed outlier pattern can be considered periodic. A more fine-grained periodicity check may involve breaking up the given day into multi-hour groups.

III. EVALUATION

1) *Forecasters*: Though ATH does not dictate the forecasting algorithm to be used to obtain the residuals, the effectiveness of time series AD depends on the accuracy of the underlying forecasting method. Quartile Based Seasonality Decomposition (QBSD) (filed as a patent by Ericsson [15]) is a live data forecasting algorithm tailored for telecom data; it does not require an explicit retraining step. QBSD has been extensively evaluated in [16] with multiple state-of-the-art forecasters. Since the focus of this paper is not to compare forecasters and detectors but to illustrate the applicability of ATH, only QBSD and Prophet [17] are applied for forecasting.

2) *Outlier Detectors*: The PyOD toolbox [18] is a Python library of popular outlier detection algorithms. 11 recent detectors in PyOD were evaluated with ATH as the thresholding technique. 6 empirically best performing models were found to be Lightweight On-line Detector of Anomalies (LODA) [19], Deep One-Class Classification (DeepSVDD) [20], Copula-Based Outlier Detection (COPOD) algorithm [21], Empirical Cumulative Distribution Functions (ECOD) [22], AutoEncoder [23] and Variational AutoEncoder (VAE) [24]. Only the results of these 6 models will be included in this paper.

The Z-Score is a statistical metric utilized to gauge the deviation of a data point from the mean of its corresponding distribution, expressed in terms of standard deviations. The Z-Score, Z , is given by $Z = (x - \mu) / \sigma$, where x is the data point, μ is the mean of the distribution, and σ is the standard deviation of the distribution.

The POT [12] method does not provide a concrete anomaly threshold; it only provides candidate anomalies, i.e., the peaks. ATH can consider each of these peaks as candidate thresholds to provide a better threshold that satisfies the heuristic. The

Python implementation of the POT algorithm [25] is utilized for this study.

3) *Candidate Datasets*: For a univariate anomaly detection problem, the observed variable should possess distinct characteristics that visually differentiates normal and anomalous points when plotted. The definition of an anomaly, according to the use case of interest in this paper, is a rare occurrence that deviates from the typical distribution of the data characterised by an abnormally high peak or abnormally low dip in value. Multiple experiments were conducted utilizing existing datasets such as the AIOps [26] and the NAB datasets [27]. However, the anomalies labelled in these datasets were found to be inconsistent with the definition of an anomaly as per the scope of this paper. Examples include labelling neighbouring non-anomalous points of an anomalous peak, or labelling a substantial proportion of the input data labelled as anomalies. Consequently, these datasets are not part of the final analysis.

4) *Selected Dataset*: The dataset utilized for the purpose of the evaluation is referred to as EON1-Cell-U, which is part of the Ericsson Outlier Nexus (EON) [28]. EON1-Cell-U is composed of time series KPI designed for univariate AD. There are 10 KPIs in this dataset encompassing different time series characteristics, both seasonal and stochastic. The first 6 KPIs (A through F) exhibit seasonality, that is, a predictable component that is dependent on the time of the day and the day of the week. This pattern can be observed in load KPIs, e.g., Active Uplink Users. The last 4 KPIs (G through J) exhibit stochasticity without seasonality which is characterized by erratic behavior. Such a pattern can be observed by fault monitoring KPIs, e.g., S1 Setup Failure. The KPIs have separate periods for training, validation, and testing; one month each. The interval between two consecutive KPI values is 15 minutes. Training split includes the month of February 2023 (28 days). Similarly, the validation and the test splits include March 2023 (31 days) and April 2023 (30 days) respectively. Anomalies for each KPI have been labeled as either 1, -1, or 0, depending on whether the occurrence is a right-tailed anomaly, a left-tailed anomaly, or not an anomaly respectively. However, these labels are not fed into the algorithm while training or testing (the thresholding is based on unsupervised learning). The labels are used to derive the evaluation metrics from the predictions of the ATH algorithm.

TABLE I
ATH PERFORMANCE ON EON1-CELL-U IN TERMS OF F_1 SCORES

Forecaster	Detector	A	B	C	D	E	F	G	H	I	J	Mean
Prophet	DeepSVDD	0.571	0.785	0.725	0.748	0.200	0.290	0.074	1.000	0.043	1.000	0.544
Prophet	LODA	0.056	0.780	0.779	0.818	0.042	0.415	0.769	0.774	0.930	0.766	0.613
Prophet	ECOD	0.556	0.484	0.773	0.608	0.328	0.296	0.600	0.923	0.653	0.909	0.613
Prophet	COPOD	0.163	0.984	0.424	0.250	0.273	0.436	0.889	0.923	0.952	0.837	0.613
Prophet	AutoEncoder	0.467	0.608	0.697	0.696	0.328	0.471	0.741	0.750	0.930	0.766	0.645
Prophet	VAE	0.467	0.608	0.697	0.696	0.328	0.471	0.741	0.750	0.930	0.766	0.645
Prophet	Z-Score	0.485	0.933	0.767	0.727	0.319	0.391	1.000	0.909	1.000	0.971	0.750
QBSD	COPOD	0.227	0.951	0.444	0.290	0.333	0.481	0.870	0.923	0.930	0.837	0.629
QBSD	ECOD	0.600	0.593	0.767	0.622	0.657	0.302	0.818	0.566	0.833	0.837	0.659
QBSD	DeepSVDD	0.625	0.821	0.840	0.870	0.732	0.240	0.833	0.727	0.976	0.621	0.729
QBSD	LODA	0.625	0.941	0.757	0.870	0.714	0.279	0.741	0.787	0.930	0.766	0.741
QBSD	AutoEncoder	0.654	0.889	0.785	0.844	0.696	0.483	0.714	0.774	0.851	0.766	0.746
QBSD	VAE	0.654	0.889	0.785	0.844	0.696	0.483	0.714	0.774	0.851	0.766	0.746
QBSD	Z-Score	0.636	0.901	0.861	0.588	0.711	0.857	0.800	0.909	0.800	0.971	0.803

The best-performing detector for each dataset has been typeset in boldface for reference.

5) *Performance Metric*: F_1 score, that is, the harmonic mean of precision and recall, is used as the assessment metric. Though these metrics are usually employed to evaluate supervised learning methods, the availability of the labelled dataset in this problem enables these metrics to be used for this use case.

IV. EXPERIMENTAL RESULTS AND DISCUSSION

The values for the *periodicity_limit* and *proportion_limit* were set based on optimal scores on the validation set. The typical value for *periodicity_limit* ranged between 2 to 4, while the *proportion_limit* ranged between 0.005 to 0.01 (i.e. less than 1% data to be anomalous). Fig. 3 depicts the effectiveness of ATH in detecting anomalies with QBSD as the forecaster and Z-Score as the detector with ground truth labels (GT). For seasonal KPIs (A-F), it would be technically possible to apply thresholding directly on the thresholds for the use case concerned. However, stochastic KPIs (G-J) require a deeper focus since these KPIs cannot be forecasted. A standard “three-sigma-rule” will mark all peaks of a such KPI as outliers even if it is not anomalous. While the periodicity condition addresses seasonal aspects that can be overlooked by forecasters, the proportionality condition addresses the stochasticity ensuring minimal false positives and false negatives.

The performance metrics computed from the predictions generated by ATH on each forecaster-detector combination for each KPI, were grouped and presented in Table I. The showcased outcomes offer valuable insights into the efficacy of ATH across different forecasters and detectors. DeepSVDD performs remarkably well in some stochastic scenarios (H and J), and some seasonal scenarios (D and E) when coupled with QBSD. Both Autoencoders (AutoEncoder and VAE) gave the same results for all scenarios. There is no single forecaster-detector combination can be considered the best for all KPIs. Nevertheless, on average, QBSD forecaster consistently outperforms the results obtained by using Prophet forecaster with the ATH algorithm regardless of the the detector used (Fig. 4).

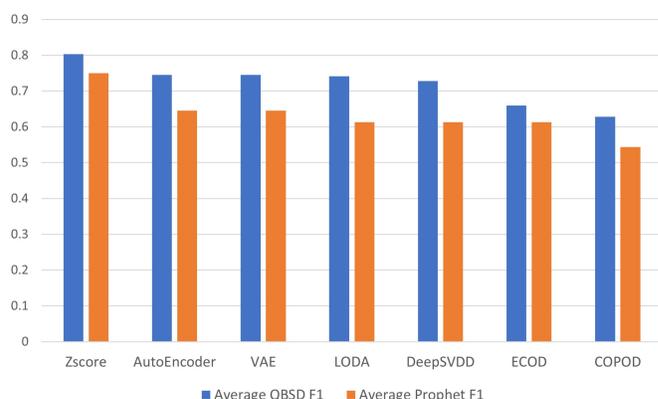


Fig. 4. Average F_1 scores across all KPIs of EON1-Cell-U

ECOD and COPOD seemed to be the least accurate detectors considered. Surprisingly, the simple Z-Score seems to perform better on average than the other specialized PyOD detectors. This phenomenon can be due to the way the detectors are designed. The outlier score as represented in the Z-Score has a scalar dependency on the magnitude of the anomalous spike. This dependency aligns well with the definition of an anomaly in the dataset considered and the working of ATH. The outlier scores produced by the PyOD algorithms considered in this experimentation does not always follow this linear dependency.

An effective AD solution should also be efficient to make it viable for business, especially in the context of big data applications. To derive the computational complexity of ATH, let n be the number of elements in the input window. Sorting is done in $O(n)$ time using Quicksort. Since the anomaly proportion limit is obviously far less than $n/2$, only $O(\log n)$ scores are candidates for thresholds. Each of the n elements are processed for every threshold candidate. Temporal difference between each outlier is computed in quadratic time to account for the periodicity check. Thus the overall computational complexity of ATH can be considered to be $O(n \log^2 n)$.

V. CONCLUSION

The proposed ATH algorithm offers a computationally efficient solution to time series AD. It incorporates business logic to differentiate between outliers and use-case-specific anomalies. Its key advantage is its ability to accommodate different outlier detectors and residual extraction methods, making it compatible with a wide range of AD algorithms. The experimental results prove its practical significance in telecom KPIs. The next step would be to extend this heuristic to capture more complex KPI interactions for special cases of anomalies.

REFERENCES

- [1] Z. Ming, C. Sun, X. Li, Q. Fan, X. Wang, and V. C. Leung, "Ensemble learning based sleeping cell detection in cloud radio access networks," in *2020 IEEE Symposium on Computers and Communications (ISCC)*. IEEE, 2020, pp. 1–6.
- [2] J. Hong, S. Park, J.-H. Yoo, and J. W.-K. Hong, "Machine learning based SLA-aware VNF anomaly detection for virtual network management," in *2020 16th International Conference on Network and Service Management (CNSM)*. IEEE, 2020, pp. 1–7.
- [3] A. S. Alghawli, "Complex methods detect anomalies in real time based on time series analysis," *Alexandria Engineering Journal*, vol. 61, no. 1, pp. 549–561, 2022.
- [4] C. Zhang, D. Song, Y. Chen, X. Feng, C. Lumezanu, W. Cheng, J. Ni, B. Zong, H. Chen, and N. V. Chawla, "A deep neural network for unsupervised anomaly detection and diagnosis in multivariate time series data," in *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, ser. AAAI'19/IAAI'19/EAAI'19. AAAI Press, 2019. [Online]. Available: <https://doi.org/10.1609/aaai.v33i01.33011409>
- [5] J.-R. Jiang, J.-B. Kao, and Y.-L. Li, "Semi-supervised time series anomaly detection based on statistics and deep learning," *Applied Sciences*, vol. 11, no. 15, 2021. [Online]. Available: <https://www.mdpi.com/2076-3417/11/15/6698>
- [6] S. Chatterjee, R. Bopardikar, M. Guerard, U. Thakore, and X. Jiang, "Mospat: Automl based model selection and parameter tuning for time series anomaly detection," *arXiv preprint arXiv:2205.11755*, 2022.
- [7] J. Audibert, P. Michiardi, F. Guyard, S. Marti, and M. A. Zuluaga, "Usad: Unsupervised anomaly detection on multivariate time series," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 3395–3404.
- [8] S. Han, X. Hu, H. Huang, M. Jiang, and Y. Zhao, "Adbench: Anomaly detection benchmark," *Advances in Neural Information Processing Systems*, vol. 35, pp. 32 142–32 159, 2022.
- [9] N. Laptev, S. Amizadeh, and I. Flint, "Generic and scalable framework for automated time-series anomaly detection," in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 1939–1947. [Online]. Available: <https://doi.org/10.1145/2783258.2788611>
- [10] N. Bezak, M. Brilly, and M. Šraj, "Comparison between the peaks-over-threshold method and the annual maximum method for flood frequency analysis," *Hydrological Sciences Journal*, vol. 59, no. 5, pp. 959–977, 2014.
- [11] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Soderstrom, "Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ser. KDD '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 387–395. [Online]. Available: <https://doi.org/10.1145/3219819.3219845>
- [12] A. Siffer, P.-A. Fouque, A. Termier, and C. Largouet, "Anomaly detection in streams with extreme value theory," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017, pp. 1067–1075.
- [13] S. Tuli, G. Casale, and N. R. Jennings, "Tranad: Deep transformer networks for anomaly detection in multivariate time series data," *Proc. VLDB Endow.*, vol. 15, no. 6, p. 1201–1214, feb 2022. [Online]. Available: <https://doi.org/10.14778/3514061.3514067>
- [14] E. R. H. P. Isaac, "Adaptive thresholding heuristic for anomaly detection," September 2021, Patent No. WO2021176460A1. [Online]. Available: <https://patents.google.com/patent/WO2021176460A1/>
- [15] E. R. H. P. Isaac, P. Bhargava, and M. Gottumukkala, "First node and methods performed thereby for handling anomalous values," June 2021, Patent No. WO2022271057A1. [Online]. Available: <https://patents.google.com/patent/WO2022271057A1/>
- [16] E. R. H. P. Isaac and B. Singh, "QBSD: quartile-based seasonality decomposition for cost-effective time series forecasting," 2023, arXiv preprint. [Online]. Available: <https://arxiv.org/abs/2306.05989>
- [17] S. J. Taylor and B. Letham, "Forecasting at scale," *The American Statistician*, vol. 72, no. 1, pp. 37–45, 2018.
- [18] Y. Zhao, Z. Nasrullah, and Z. Li, "PyOD: A python toolbox for scalable outlier detection," *arXiv preprint arXiv:1901.01588*, 2019.
- [19] T. Pevný, "Loda: Lightweight on-line detector of anomalies," *Machine Learning*, vol. 102, pp. 275–304, 2016.
- [20] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft, "Deep one-class classification," in *International conference on machine learning*. PMLR, 2018, pp. 4393–4402.
- [21] Z. Li, Y. Zhao, N. Botta, C. Ionescu, and X. Hu, "Copod: copula-based outlier detection," in *2020 IEEE international conference on data mining (ICDM)*. IEEE, 2020, pp. 1118–1123.
- [22] Z. Li, Y. Zhao, X. Hu, N. Botta, C. Ionescu, and G. Chen, "Ecod: Unsupervised outlier detection using empirical cumulative distribution functions," *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- [23] C. C. Aggarwal, "Outlier analysis," in *Data Mining*. Springer, 2015, ch. 3, pp. 75–79.
- [24] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [25] G. Bocharov, "pyextremes – Extreme Value Analysis (EVA) in Python," 2022. [Online]. Available: <https://georgebv.github.io/pyextremes/>
- [26] Z. Li, N. Zhao, S. Zhang, Y. Sun, P. Chen, X. Wen, M. Ma, and D. Pei, "Constructing large-scale real-world benchmark datasets for aiops," *arXiv preprint arXiv:2208.03938*, 2022.
- [27] A. Lavin and S. Ahmad, "Evaluating real-time anomaly detection algorithms—the numenta anomaly benchmark," in *2015 IEEE 14th international conference on machine learning and applications (ICMLA)*. IEEE, 2015, pp. 38–44.
- [28] Ericsson Research, "Ericsson Outlier Nexus (EON)," <https://github.com/EricssonResearch/eon>, accessed: 2023-06-27.