

State-of-the-Art Report

Anomaly detection using Autoencoders and Methods to Improve them

Cozmina Scorobete

Advisory: Darian Onchiş

Facultatea de Matematică şi Informatică
Universitatea de Vest, Timişoara, România

Abstract

The rapid growth of the technology made it crucial to have systems protected against cyber-attacks. Machine Learning (ML) and Deep Learning (DL) techniques have proven to be effective in network anomaly detection. Explainable Artificial Intelligence (XAI) offers an approach in which we can understand those "black boxes", furthermore addressing the main limitation of ML and DL, the lack of trust in the systems. This SoTA focuses on a specified paper [RZ21] that explores the use of Shapley Additive Explanation (SHAP), an XAI technique, to improve the transparency and effectiveness of autoencoder-based models for network anomaly detection.

Contents

1	Introduction	3
1.1	Motivation	3
1.2	Context of the Research	3
1.3	Aim of the Report	4
1.4	Structure of the Report	5
2	Methodology	5
2.1	Selection of Bibliographical Sources	5
2.2	Statistical Overview	6
3	State of the Art	8
3.1	Papers-Set 1 : Reviews on Autoencoders	8
3.2	Papers-Set 2 : XAI	8
3.3	Papers-Set 3 : SHAP	9
3.4	Papers-Set 1 : Reviews on Autoencoders	9
3.5	Papers-Set 2 : XAI	9
3.6	Papers-Set 3 : SHAP	9
3.7	Papers-Set 3 : Anomaly Detection Using Autoencoders	10
3.8	Anomaly Detection with Robust Deep Autoencoders	10
3.9	A study of deep convolutional auto-encoders for anomaly de- tection in videos	11
3.10	Anomaly Detection Using Autoencoders with Nonlinear Di- mensionality Reduction	11
3.11	Unsupervised Anomaly Detection in Time Series Using LSTM- Based Autoencoders	12
3.12	Utilizing XAI Technique to Improve Autoencoder Based Model for Computer Network Anomaly Detection with Shapley Ad- ditive Explanation(SHAP)	12
4	Critical Analysis	12
4.1	Comparative Analysis	13
4.2	Performance analysis on my model vs the one in [RZ21]	14
5	Conclusions and Further Work	16
5.1	Summary	16
5.2	Future Directions	16

1 Introduction

Due to the increasing dependence on digital systems and internet connectivity, there has been a growth in malicious attacks. As a result, anomaly detection has become a critical field in cybersecurity. The aim of anomaly detection is to identify unusual traffic patterns that indicate potential threats, such as Distributed Denial-of-Service (DDoS), infiltration and data exfiltration. However, while machine learning (ML) and deep learning (DL) methods have shown significant promise in detecting anomalies, they are often criticized for their lack of interpretability. These models, capable of learning from unlabeled data, offer unique opportunities for identifying subtle and emerging patterns in network traffic that traditional systems might overlook.

1.1 Motivation

This report was inspired by my personal interest in machine learning and its applications in anomaly detection, which I recently discovered while working on another project involving anomaly detection on the stock market using unsupervised learning, particularly the possibilities of unsupervised learning techniques such as autoencoders. I noticed two critical challenges during my previous work: the inability to clearly explain why certain anomalies were flagged and the limitations in optimizing model performance.

These challenges are not unique to the stock market but also in other domains like network security. The method used in [RZ21] came with the realization that explainability could not only build trust in the model but also improve the performance, which motivated me to look deeper in the intersection of Explainable Artificial Intelligence and autoencoder-based anomaly detection challenges.

1.2 Context of the Research

The field of anomaly detection presents several questions that I aim to address, however some of them still remain to be worked on in future research topics.

1. How can machine learning models detect anomalies in complex systems (network traffic)?
2. What methods can be used to add interpretability to these models ensuring that the decisions are transparent?
3. Can XAI techniques such as Shapley Additive Explanation (SHAP) improve the accuracy of autoencoders or other models?

4. How can unsupervised learning models address the lack of labelled data while maintaining high detection performance?

Some real world problems are DDos Attacks, Zero-Day Vulnerabilities, Web Exploits and Sql Injections.

- DDos Attacks: Distributed Denial-of-Service (DDoS) attacks flood network servers with malicious traffic, making it difficult to distinguish between legitimate and abnormal activities. This kinds of attacks appear in the data used to train the model used in the [RZ21], as well as in my own model that tried to mimic the work of the authors to see if indeed the SHAP is an improvement on unsupervised learning based on autoencoders.
- Web Exploits and SQL Injections: These targeted attacks can disrupt services or compromise sensitive data, requiring robust detection systems capable of real-time response.

1.3 Aim of the Report

The aim of this report is to have a broad analysis on the field of anomaly detection focusing on the combination between XAI and traditional anomaly detection. A comprehensive analysis on the current state of the art method and their applications.

1. Critically Analyze the Field: Evaluate existing approaches and tring to replicate the work specified in [RZ21]. The focus is on unsupervised learning techniques like autoencoders.
2. Compare Different Approaches: In order to better understand if XAI methods may enhance the performance of the model two approaches are taken into account the traditional methods and those enhanced by XAP
3. Highlight Open Problems: To identify gaps in current research, including challenges that may have been overlooked by the researchers in their papers.
4. Provide a Foundation for Future Research: Based on the analysis, propose potential directions for future work on improving the ML models using XAI methods.

1.4 Structure of the Report

1. Introduction: In this section I provided the motivation behind the report, the context of the research the the aim of the report as well as the structure of the report.
2. Methodology: This section provides on overview on how the papers were selected and analyzed. As well as the selection of the datasets. It also includes a selection criteria such as sources, relevance and quality.
3. State of the art in the field: In this section are provided summaries of the selected papers and and overview of the filed as it is today.
4. Critical Analysis: This section compares different approaches presented in the selected papers as well an comparison with my own results from the method presented in [RZ21]. The evaluation is based on accuracy. It highlights the limitations and the advantages of each approach.
5. Conclusion and further work: In the final section are presented the most relevant findings and offers recommendations for future research.

2 Methodology

2.1 Selection of Bibliographical Sources

In order to have a solid foundation for the report, a wide range of sources were consulted. And each papers needed to meet a number of criteria in order to be included in the report.

Some sources are: Science Direct, Google Scholar,DBLP, arXiv, Research Gate. The key words used were:


- Anomaly Detection: to focus on the hole field
- Autoencoders: in order to have a better understanding on how they work
- Explainable Artificial Intelligence(XAI)
- Shapley Additive Explanation (SHAP): to better understand the technique

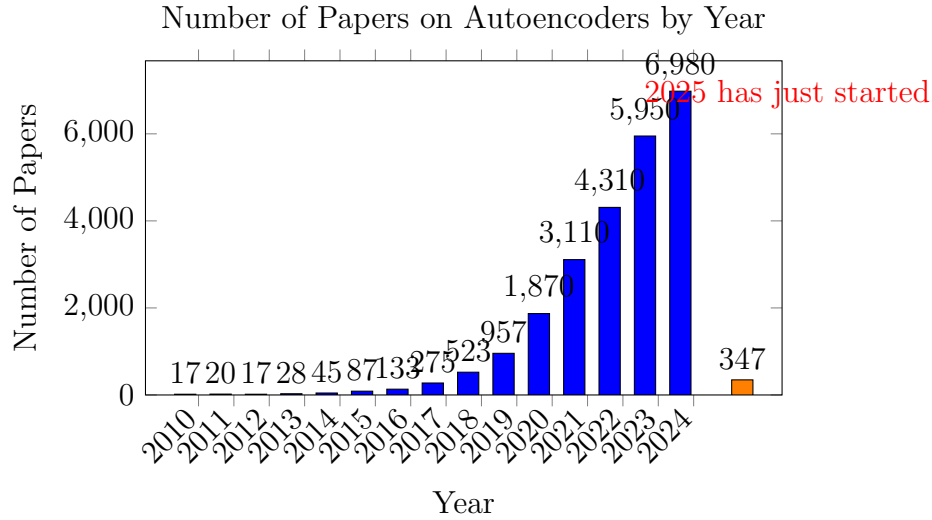
Some inclusion Criteria are: papers published by reputable journals, research that present practical applications not only theoretical, publishing


year to be at least in the 10-15 years, the number of citations must be higher then 50.

Some exclusion criteria are: papers lacking methodology or papers that have questionable results, studies that are not relevant to the topic of the report and outdated works.



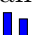
2.2 Statistical Overview

In order to see the interest in the field of autoencoders I conducted an analysis focused on the number of papers that was published during the years 2010-2025 the analysis can be seen in the fallowing bar chart .

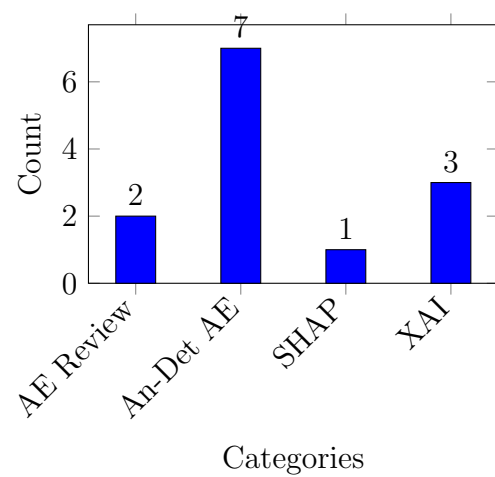


In  we can clearly see that the interest in this field has a massive growth in the last 2 years, not only that but in the first 17 days of 2025 there are already published 347 on the topic of autoencoders(AE).

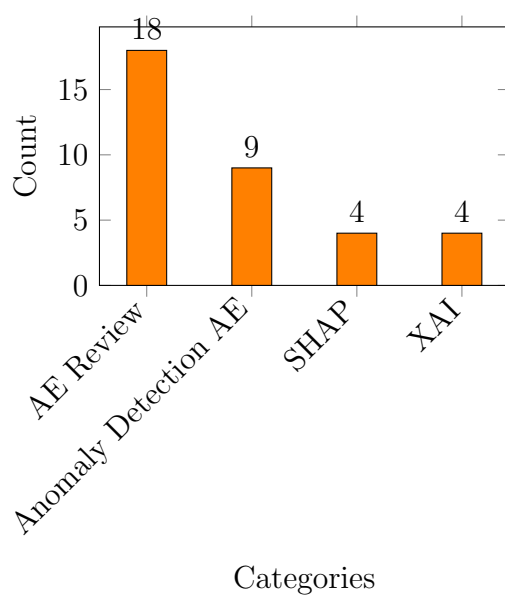
After the overview of the topic presented above a more in depth selection was conducted, in which for each key word I applied the acceptance criteria presented above and I selected a small number of papers that meet the criteria as seen in the plots below:

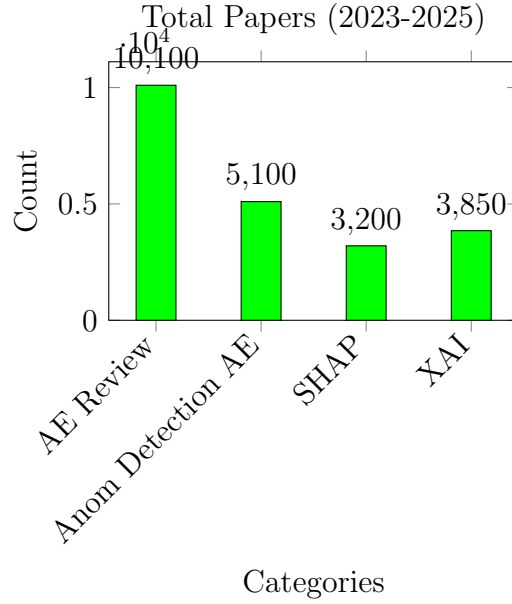
In  can be seen the number of papers that appeared on google scholar after each key word, after that I have decided to select a number of maximum 20 paper as seen in  that meet the criteria of being published on a trustworthy source and number of citation at least 50. The last round of selection as seen in  was based on the main topic of each paper and how relevant it would be for this study.

Selected Papers (2023-2025, 2020-SHAP)



Papers Meeting Criteria (2023-2025)





3 State of the Art

3.1 Papers-Set 1 : Reviews on Autoencoders

In this subsection we will discuss the first set of papers that are [BDS⁺24] and [CG23]. The selected papers focus on Auto encoders and their variants, describing the basic architecture of such a model as well as their variants. They highlight the ability of autoencoders to deal with noise or data compression. In both papers are also mentioned challenges such as overfitting, generalization and sensitivity to hyperparameters.

3.2 Papers-Set 2 : XAI

In this subsection we will discuss the field of XAI as presented in this papers [CHS⁺23], [CSTT24] and [CKAA23]. The main topic presented in those papers is how can we make trustworthy and transparent algorithms for all users with focus on their applications in practical fields. All papers present methods that take into account the human factor as well as different approaches of model-agnostic methods, white-box models and post-hoc explanations. The main results are the practical examples of XAI, the challenges presented and the disproportionality between interpretability and accuracy. The open problems mentioned are not unknown to this field they include challenges such as scalability, bias control and the human component.

3.3 Papers-Set 3 : SHAP

3.4 Papers-Set 1 : Reviews on Autoencoders

In this subsection we will discuss the first set of papers that are [BDS⁺24] and [CG23]. The selected papers focus on Auto encoders and their variants, describing the basic architecture of such a model as well as their variants. They highlight the ability of autoencoders to deal with noise or data compression. In both papers are also mentioned challenges such as overfitting, generalization and sensitivity to hyperparameters.

3.5 Papers-Set 2 : XAI

In this subsection we will discuss the field of XAI as presented in this papers [CHS⁺23], [CSTT24] and [CKAA23]. The main topic presented in those papers is how can we make trustworthy and transparent algorithms for all users with focus on their applications in practical fields. All papers present methods that take into account the human factor as well as different approaches of model-agnostic methods, white-box models and post-hoc explanations. The main results are the practical examples of XAI, the challenges presented and the disproportionality between interpretability and accuracy. The open problems mentioned are not unknown to this field they include challenges such as scalability, bias control and the human component.

3.6 Papers-Set 3 : SHAP

This section focuses on the paper on SHAP [SHJ⁺20]. This paper presents the vulnerabilities in SHAP and Lime. In the paper are presented multiple methods to 'trick' the methods(LIME, SHAP). One of them is an adversarial classifier it changes the explanation to focus on fake features rather than on real ones like race or gender, another approach is Out-of-Distribution Detector this is a detector for real data or perturbations. The results shown that SHAP as well as LIME was fooled and it presents that LIME it is more likely to be vulnerable to this attacks than SHAP but the second one also fails some tests. Some open problems are generalization of the stack to study more methods and to implement some defensive mechanisms against such stacks.

3.7 Papers-Set 3 : Anomaly Detection Using Autoencoders

The last set of papers is the biggest one. It focuses on articles that present anomaly detection in different fields. We will discuss here the main ideas of those papers and in section 4 I will present an in depth analysis of some of the methods presented in those papers.

All papers address similar problems concerning the detection of anomalies in high-dimensional, noisy or temporal datasets using autoencoders that proved useful in such settings where the labels are unavailable.

Some of the methods that were most effective are LSTM-based (Long Short Term Memory) Autoencoders, Sparse Autoencoders, Convolutional Autoencoders and Attention Mechanisms, there were also mentioned methods like combining autoencoders with other models like Gaussian Mixture Models.

Some of the common results in those papers were the high accuracy presented and a robust detection. There are also a lot of challenges to remain unsolved as this is a field that grew in popularity in the recent years, some of them are:

- Scalability: it is computationally intensive on large sequences or datasets
- Parameter Sensitivity those kind of models are very sensitive to their hyperparameters like number of layers or step sizes

The papers we have discussed in this subsection are: [ZDWP⁺24], [ZP17a], [TBJS20], [SY14], [BBRSR23], [DLH23] and [PLV19]

3.8 Anomaly Detection with Robust Deep Autoencoders

In their paper [ZP17b], to overcome the typical problem of handling noisy and outlier-infested data, Zhou and Paenroth provide a unique approach for identifying anomalies that makes use of deep learning. In the presence of abnormalities, traditional deep-denoising autoencoders perform much worse. However, they are quite good at learning representations from clean data. The authors suggest the Robust Deep Autoencoder (RDA) to overcome this restriction, which makes conventional autoencoders more resilient.

The authors further extend their approach by introducing the Group Robust Deep Autoencoder (GRDA), in order to handle more complicated situations in which the data include both random anomalies and organized corruption. The advantage of this approach is that it can distinguish between

the types of noise. This would be helpful in real world applications where the noise may not be random.

The performance of the model is proven by extensive experiments on benchmarks datasets. They evaluated it in contrast with other state-of-the-art anomaly detection methods, and the results show that their approach performs better in terms of accuracy and robustness.

3.9 A study of deep convolutional auto-encoders for anomaly detection in videos

Perlin and Lopes [RLL18] address the significant limitations of traditional handcrafted features used in anomaly detection for video surveillance. Methods like social force models, dense trajectories, and dynamic textures require a broad domain-specific knowledge, which is hard to use on different applications. This reliance on a prior knowledge often leads to poor performance. The problem with those traditional approaches is the need of predefined features and assumptions, leading to an approach that is automatic and can generalize features from data. They propose the use of Convolutional Auto-Encoders (CAEs) in order to identify anomalies in videos. They have used reconstruction error as anomalies, but also tested how more information like motion data added to the frames affected the model's performance. A series of experiments were conducted and the results show that if appearance features are added to the models then the performance increases.

3.10 Anomaly Detection Using Autoencoders with Non-linear Dimensionality Reduction

In this paper [SY14], appears the idea of detecting anomalies in datasets with nonlinear relationships using autoencoders. They use them to reduce the dimension of the data while still maintain complex relations. They tested the system on synthetic data (Lorenz System) as well as on real world data from spacecraft telemetry data. They have compared the performance of the model with the linear PCA for detecting anomalies and the proposed model outperformed linear PCA. This study also indicates the efficiency of the autoencoder compared to that of the PCA. By analyzing all they demonstrate the efficiency of the model to characterize the anomalies.

3.11 Unsupervised Anomaly Detection in Time Series Using LSTM-Based Autoencoders

This paper [PLV19] investigates unsupervised anomaly detection in time series using Long Short-Term Memory (LSTM)-based autoencoders. The difficulty in this field is challenging due to the unpredictable nature of anomalies and the lack of labeled datasets. The traditional methods (clustering or Gaussian distribution-based techniques) do not seem to categorize all the dependencies.

The authors propose using autoencoders in order to reduce dimension. The study focuses on the advantage of LSTM layers in keeping temporal information. The method was tested on artificial datasets as well as on DCASE dataset and it shown an accuracy of 87%, but only after applying smoothies and response adjustment . The accuracy shows how hard is to optimize a model, the biggest challenge being the selection of hyper-parameters.

3.12 Utilizing XAI Technique to Improve Autoencoder Based Model for Computer Network Anomaly Detection with Shapley Additive Explanation(SHAP)

In this paper[RZ21] the objective is to prove that XAI techniques can be used to improve the performance of an autoencoder-based anomaly detection model. The proposed approach is to compare 3 different models, one traditional trained on all features, one that applies feature selection and a third one that is build on the top 30 features identified using SHAP values. The result shoes that the SHAP model has a higher accuracy and minimizes irrelevant features. To calculate this it uses ROC-AUC metric. It also states the challenge of this approach which is the computational cost that is very high due to the kernel SHAP method. The open question is to apply SHAP to other unsupervised deep learning architectures.

4 Critical Analysis

I would like to analyses the strengths and challenges of some models presented in the papaers discussed in subsection 3.4. We will compare LSTM-Based Autoencoders as presented in ?? this is a method that uses LSTM layers to learn temporal dependencies, ideal for tasks like anomaly in network traffic. The second method we will analyse is a standard autoencoder presented in [SY14] this method uses nonlinear dimensionality reduction with

autoencoders. The third approach we will discuss is Robust Deep Autoencoder(RDA) from [ZP17a] combines the autoencoders with robust PCA to decompose the output in low-rank components. And the 3 methods presented in [RZ21], first approach trains on all features, the second approach applies feature selection, shap approach uses shap for feature selection. For a detailed information see 1

4.1 Comparative Analysis

Method	Challenges	Strengths	Metrics
Standard Autoencoder [SY14]	1. Struggles in datasets with high noise. 2. Does not have feature correlation	1. Good on non-linear datasets. 2. Efficient from a computation view.	Lorenz Dataset: Autoencoder AUC = 0.85 vs PCA AUC = 0.6. It improved the anomaly detection with PCA
RDA [ZP17a]	1.High computational cost. 2.Sensitive to the hyperparameters	1. Effective in noisy data.	MNIST Digits: 30% better anomaly detection accuracy than standard autoencoders for high noise data. Structured Anomalies: F1-score = 0.64 compared to Isolation Forest (F1 = 0.37).
LSTM Autoencoder [PLV19]	1. Performance correlated to hyperparameter tuning .	1. Effective on sequential datasets. 2.Has a better performance then classic methods	Rare Sound Detection: 87% accuracy..

Method	Challenges	Strengths	Metrics
SHAP_Model [RZ21]	1. High computational cost 2. Requires careful selection of dataset.	1. Improves anomaly detection by focusing on key features. 2. Provides interpretability by explaining which features contribute to anomalies.	CICIDS2017 Dataset: Accuracy = 94%, AUC = 0.969
Model_1 [RZ21]	1. Inefficient due to redundant and irrelevant features. 2. Lower performance .	1. Can handle all features without preprocessing or selection.	CICIDS2017 Dataset: Accuracy = 81.9%, AUC = 0.819.
Model_2 [RZ21]	1. May exclude important features. 2. Redundant features might not always degrade performance, leading to inconsistent results.	1. Reduces feature space, improving model efficiency.	CICIDS2017 Dataset: Accuracy = 84.3%, AUC = 0.843.

Table 1: Comparison between approaches

As seen in table 1 the SHAP models seems to be a good approach to add interpretability as well as grow the accuracy, however this comes with the price of being computationally demanding. The RDA was the best method for separating anomalies from noise where the LSTM methods is good in situations that include time series dependencies.

4.2 Performance analysis on my model vs the one in [RZ21]

In 2 we can see clearly that the SHAP model implemented by me had a worse accuracy. However this is not conclusive due to the fact that I was not able to use the full dataset due to hardware. The SHAP requires a bigger computation power than the one that I had available, this fact only making it more clear that even tho SHAP may be a good way of model improvement it comes with great costs.

However in the table I did not presented the Precision, Recall and F1 scores for the first approach were all 0 clearly indicating the bias towards normal data. For the SHAP model the Precision is 0.1529, Recall is 0.6842 and the F1 is 0.2500 which is not a great case but better than the one seen before

Method	Paper	My Result
SHAP_Model (2021)	Accuracy = 94%, AUC = 0.969.	, AUC = 0.7194.
SHAP_Model (Baseline): Model_1	Accuracy = 81.9%, AUC = 0.819.	Accuracy = 93.17%, AUC = 0.9740..
Model_2	Accuracy = 84.3%, AUC = 0.843.	Not evaluated.

Table 2: The result from training on a part of the dataset

The following table 3 presents the models performance on the full dataset. WE can clearly see that the performance is lower this could be because of the hardware, the full training having a dataset of 2GB. In this case for both models the metrics Precision, Recall and F1 score were also 0. This is still not conclusive, I believe the model would have a better performance on better hardware.

Method	Paper	My Result
SHAP_Model (2021)	Accuracy = 94%, AUC = 0.969.	, Could not compute.
SHAP_Model (Baseline): Model_1	Accuracy = 81.9%, AUC = 0.819.	Accuracy = 42%, AUC = 0.65568
Model_2	Accuracy = 84.3%, AUC = 0.843.	Accuracy = 42%, AUC = 0.5867.

Table 3: The result from training on the full dataset

As seen in tables 3 and 2 the model that I have implemented is clearly biased duo to the fac that it is only trained on normal data. However the models presented in the paper do not seem to be ass biased as my own model.

5 Conclusions and Further Work

This report proposed an analysis on the performance of autoencoders in anomaly detection as well as an review of a method to improve autoencoders using SHAP.

5.1 Summary

The comparative analysis shewed that SHAP could improve a model but that still remain to be decided in further research. I have analyzed different kind of models like LSTM-based autoencoders, RDA and Robust Deep Autoencoders, the paper present an analysis with their streamlets and challenges.

There were also some limitations provided for the SHAP model one of the biggest challenges being the computation effort required to use such a model.

5.2 Future Directions

For future research I would like to try a different hardware to test if the SHAP model could be trained with such a large dataset. I would also like to use data that contains attacks when training the model to combat the biased results.= and last but not least I would like to see how the other methods presented work with this data set or in anomaly detection in stock market.

References

- [BBRSR23] Mohammed Ayalew Belay, Sindre Stenen Blakseth, Adil Rasheed, and Pierluigi Salvo Rossi. Unsupervised anomaly detection for iot-based multivariate time series: Existing solutions, performance analysis and future directions. *Sensors*, 23(5):2844, 2023.
- [BDS⁺24] Kamal Berahmand, Fatemeh Daneshfar, Elaheh Sadat Salehi, Yuefeng Li, and Yue Xu. Autoencoders and their applications in machine learning: a survey. *Artificial Intelligence Review*, 57(2):28, 2024.
- [CG23] Shuangshuang Chen and Wei Guo. Auto-encoders in deep learning—a review with new perspectives. *Mathematics*, 11(8):1777, 2023.
- [CHS⁺23] Vinay Chamola, Vikas Hassija, A Razia Sulthana, Debshishu Ghosh, Divyansh Dhingra, and Biplab Sikdar. A review of trustworthy and explainable artificial intelligence (xai). *IEEE Access*, 2023.
- [CKAA23] Tobias Clement, Nils Kemmerzell, Mohamed Abdelaal, and Michael Amberg. Xair: A systematic metareview of explainable ai (xai) aligned to the software development process. *Machine Learning and Knowledge Extraction*, 5(1):78–108, 2023.
- [CSTT24] Ching-Hua Chuan, Ruoyu Sun, Shiyun Tian, and Wan-Hsiu Sunny Tsai. Explainable artificial intelligence (xai) for facilitating recognition of algorithmic bias: An experiment from imposed users’ perspectives. *Telematics and Informatics*, 91:102135, 2024.
- [DLH23] Huu-Thanh Duong, Viet-Tuan Le, and Vinh Truong Hoang. Deep learning-based anomaly detection in video surveillance: A survey. *Sensors*, 23(11):5024, 2023.
- [PLV19] Oleksandr I Provotar, Yaroslav M Linder, and Maksym M Veres. Unsupervised anomaly detection in time series using lstm-based autoencoders. In *2019 IEEE International Conference on Advanced Trends in Information Theory (ATIT)*, pages 513–517. IEEE, 2019.

- [RLL18] Manassés Ribeiro, André Eugênio Lazzaretti, and Heitor Silvério Lopes. A study of deep convolutional auto-encoders for anomaly detection in videos. *Pattern Recognition Letters*, 105:13–22, 2018.
- [RZ21] Khushnaseeb Roshan and Aasim Zafar. Utilizing xai technique to improve autoencoder based model for computer network anomaly detection with shapley additive explanation (shap). *arXiv preprint arXiv:2112.08442*, 2021.
- [SHJ⁺20] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186, 2020.
- [SY14] Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear dimensionality reduction. In *Proceedings of the MLSDA 2014 2nd workshop on machine learning for sensory data analysis*, pages 4–11, 2014.
- [TBJS20] Srikanth Thudumu, Philip Branch, Jiong Jin, and Jugdutt Singh. A comprehensive survey of anomaly detection techniques for high dimensional big data. *Journal of Big Data*, 7:1–30, 2020.
- [ZDWP⁺24] Zahra Zamanzadeh Darban, Geoffrey I Webb, Shirui Pan, Charu Aggarwal, and Mahsa Salehi. Deep learning for time series anomaly detection: A survey. *ACM Computing Surveys*, 57(1):1–42, 2024.
- [ZP17a] Chong Zhou and Randy C Paffenroth. Anomaly detection with robust deep autoencoders. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 665–674, 2017.
- [ZP17b] Chong Zhou and Randy C Paffenroth. Anomaly detection with robust deep autoencoders. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 665–674, 2017.