

Introducing a Multi-Perspective xAI Tool for Better Model Explainability

Marek Pawlicki^{1,2}, Damian Puchalski¹, Sebastian Szelest¹, Aleksandra Pawlicka^{1,3}, Rafał Kozik^{1,2}, and Michał Choraś^{1,2}

¹ITTI Sp. z o.o., Poznań, Poland

²Bydgoszcz University of Science and Technology, Bydgoszcz, Poland

³University of Warsaw

ABSTRACT

This paper introduces an innovative tool equipped with a multi-perspective, user-friendly dashboard designed to enhance the explainability of AI models, particularly in cybersecurity. By enabling users to select data samples and apply various xAI methods, the tool provides insightful views into the decision-making processes of AI systems. These methods offer diverse perspectives and deepen the understanding of how models derive their conclusions, thus demystifying the "black box" of AI. The tool's architecture facilitates easy integration with existing ML models, making it accessible to users regardless of their technical expertise. This approach promotes transparency and fosters trust in AI applications by aligning decision-making with domain knowledge and mitigating potential biases.

Keywords: machine learning; AI explainability; network intrusion detection

1 INTRODUCTION

In a variety of domains and applications, the accuracy of Artificial Intelligence (AI) algorithms has surpassed many traditional methods of classification and pattern recognition in data streams. This advancement is also beneficial in cyberattack detection, where, despite different attacker motivations (Pawlicka et al., 2021, 2020), AI shows great potential (Kaur et al., 2023; Rafy, 2024; Pawlicki et al., 2023; Xu et al., 2021). The majority of the high-performing AI algorithms work as black-boxes – the data is fed to the algorithm, and the algorithm provides highly accurate output, however with no indications or explanations regarding the basis of these outputs and the factors influencing the decision-making process (Choraś et al., 2020). In many applications, this knowledge is very valuable, both for AI model users, as well as model developers. Therefore, explainability of AI (xAI) and interpretation of black-box algorithms are of crucial importance for the proliferation of AI technologies in the cybersecurity domain, as one of the ways to garner user trust (Barredo Arrieta et al., 2020; Ribeiro et al., 2016; Pawlicki et al., 2024).

This paper introduces an innovative tool featuring a user-friendly dashboard designed to significantly enhance interpretability for end-users. This dashboard allows users to selectively examine data samples classified by the AI model and choose from an array of explanation methods. Each method provides a distinct perspective into the model's decision-making processes, effectively demystifying the 'black box'. This multi-perspective capability, accessible via an intuitive interface, represents a key contribution of this work, offering users a tailored and insightful exploration of AI outputs.

The rest of the paper is structured as follows:

- Section 2 is focused on the description of related work in the context of applying different methods of xAI.
- The goal of section 3 is to introduce our innovative explainer for AI-based NIDS (Network Intrusion Detection).
- Section 4 presents the local and global explainability capabilities provided by the proposed tool.
- Section 5 concludes the paper.

2 RELATED WORK

A number of recent publications show a successful application of explainable AI models used to detect or classify phenomena for cybersecurity purposes. The need for explainability capabilities is emphasized,

37 taking into account aspects of ensuring the trustworthiness, interpretability, robustness and performance of
38 AI-based detection models at the local level for each sample, and in the global perspective of the entire
39 decision-making process of the AI model (Mane and Rao, 2021; Hariharan et al., 2023; Neupane et al.,
40 2022). The importance of xAI algorithms is unquestionable for cybersecurity operators using black-box
41 models, but also crucial for the performance of the developers of AI tools for security solutions, ensuring
42 AI reliability or fine-tuning the solutions.

43 The authors of (Abou El Houda et al., 2022) presented the framework providing real-time intrusion
44 detection for Internet of Things (IoT) networks, enhanced by three xAI techniques, i.e., RuleFit, Local
45 Interpretable Model-Agnostic Explanations (LIME), and SHapley Additive exPlanations (SHAP), on top
46 of a DNN-based detection model. The aim of (Warnecke et al., 2020) was to investigate six popular
47 explanation methods and to evaluate their performance in malware detection and vulnerability discovery
48 scenarios. Assessed methods included LIME, SHAP, Gradients and Integrated Gradients (IG), Local
49 Explanation Method using Nonlinear Approximation (LEMNA) and Layer-Wise Relevance Propagation
50 (LRP). The authors investigated those methods introducing such criteria as descriptive accuracy, descriptive
51 sparsity, as well as completeness, efficiency, robustness and stability of explanation. The works conclude
52 that there exist significant differences between the xAI methods performance based on the particular
53 AI-assisted security task. On the other hand, the authors of (Arreche et al., 2024) proposed an end-to-end
54 framework to evaluate xAI methods performance in network intrusion detection tasks, including global
55 and local explanations. The authors proposed and analyzed different metrics to showcase the differences
56 between popular black-box xAI techniques, namely SHAP and LIME.

57 Additionally, several open-source, multi-method xAI frameworks, libraries and tools aiming at im-
58 proved explainability of AI models are available. XAITK toolkit is comprised of variety of separate tools
59 developed at different levels of maturity and addressing different xAI tasks, e.g.: After Action Review for
60 AI, Bayesian teaching for xAI, Counterfactual explanations, fault-line image explanations, similarity-based
61 saliency maps, etc. (Hu et al., 2021). The explAIner framework (Spinner et al., 2019), a tool for inter-
62 active and explainable machine learning, incorporates over 20 state-of-the-art xAI methods and helps to
63 understand xAI process mapping it into an iterative, three-stage workflow, including model understanding,
64 diagnosis, and refinement. OmniXAI is an open-source Python library designed to address explaining
65 decisions made by AI models to improve analytic, debugging and interpretability capabilities in a variety
66 of tasks and AI applications (Yang et al., 2022). It includes a range of well-known explanation methods,
67 such as feature-attribution/importance explanation (LIME, SHAP, Integrated Gradients, Grad-CAM, L2X),
68 counterfactual explanation (MACE), Partial Dependence Plots (PDP), and model-specific methods (linear
69 and tree models).

70 **3 EXPLAINER FOR AI-BASED NETWORK INTRUSION DETECTION DECI-** 71 **SIONS**

72 The tool proposed in this paper has been developed within the AI4CYBER (Trustworthy Artificial In-
73 telligence for Cybersecurity Reinforcement and System Resilience) project, co-funded by the European
74 Commission. The project provides next-generation trustworthy cybersecurity services that leverage AI
75 and Big Data technologies (ai4, 2024; cor, 2024). The aim of the project is to support system developers
76 and operators in effectively managing robustness, resilience, and dynamic response against advanced and
77 AI-powered cyberattacks. The solution presented in this publication is a sub-component developed within
78 the toolkit focusing on the aspects of the AI4CYBER services trustworthiness, namely explainability
79 (presented TRUST4AI.XAI explainer).

80 The high-level architecture of the solution is presented in Fig. 1. The information flow starts from
81 the user, who having authenticated, can see the dashboard and issue requests. The tool is based on the
82 publish-subscribe architecture, using the bus to move samples around the different microservices. AI
83 Model Integration Interface is implemented to handle communication with different AI models used in
84 the project, and the integration microservice is used to poll the models when the particular samples are
85 being explained. Next, different pre-processing services are used to transform the data, to tailor it to the
86 different xAI algorithms. Then, the actual xAI microservices are available through, and since some of them
87 are computationally demanding, a data storage component is implemented to save both the time and the
88 resources - when the user needs to see the same explanation multiple times.

89 The AI-based intrusion detection component, which feed the data to the explainer module, has been
90 trained on the LITNET benchmark dataset (Damasevicius et al., 2020). The xAI module allows the user to
91 view the decision making process of the AI-based detector from multiple perspectives, listed as options in
92 the dashboard.

93 In the proposed TRUST4AI.XAI component providing explainability features for AI4Cyber services
94 several state-of-the art explanation techniques, along with proprietary methods are implemented to demon-

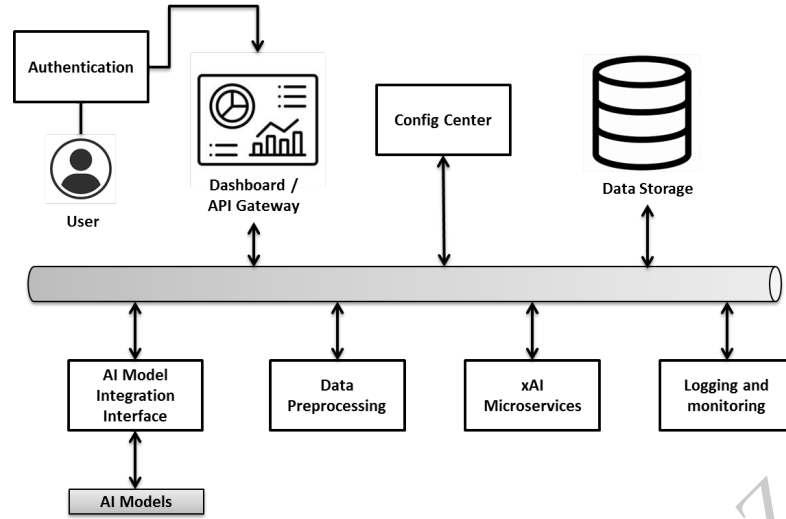


Figure 1. Architecture of the proposed xAI module.

strate different xAI capabilities that can be tailored or customized to different scenarios. There are:

1. Local explainers:

- Anchors explanations
- Diverse Counterfactual Explanations (DICE)
- Local Interpretable Model-Agnostic Explanations (LIME)
- Shapley Additive Explanations (SHAP)
- Explanation based on decision trees

2. Global explainers:

- Accumulated Local Effects (ALE)
- Individual Conditional Expectation (ICE)
- Partial Dependence Plot (PDP)
- Permutation Feature Importance (PFI)
- RuleFit method

The summary of the characteristics of each implemented explanation technique and relevant TRUST4AI.XAI dashboard views are presented in the following sub-sections.

4 LOCAL EXPLAINERS

4.1 Anchors explanations

Anchors is a model-agnostic explanation method based on if-then rules, called “anchors”. These rules are specific conditions that, if satisfied, guarantee the same prediction (model output) with a high probability (Ribeiro et al., 2018). Anchors explainer offers insights into a set of factors that are most influential in the AI decision-making process and allows for analysis of the conditions that consistently lead to the same predictions. The authors of this approach demonstrated the usefulness of anchors in different machine learning tasks, namely classification, structured prediction or text generation, for different formats of input data (i.e. tabular, text, images). According to (Ribeiro et al., 2018), the extraction of the partial input data that is sufficient for the classifier to make the prediction, helps to analyze the output of the prediction in an intuitive manner. An example of the AI explanation using this method in the proposed tool is presented in Fig. 2.

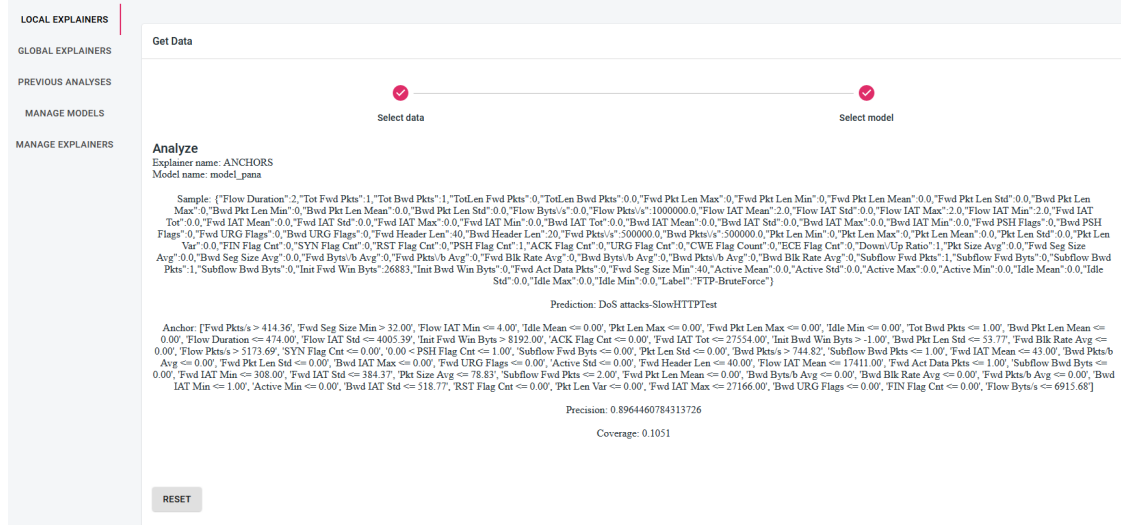


Figure 2. Anchors explanation example.

4.2 Diverse Counterfactual Explanations (DICE)

DICE (Diverse Counterfactual Explanations) is a method based on the generation of counterfactuals – hypothetical scenarios showing what minimal changes are needed to obtain different predictions. In other words, this method provides a set of different “if-then” predictions by perturbing input information, to understand the ML model’s decision process and its boundaries (Mohtilal et al., 2020). In the proposed explainer module, the user can analyze the original outcome of the model and several feasible counterfactual examples, as can be seen in Fig. 3.

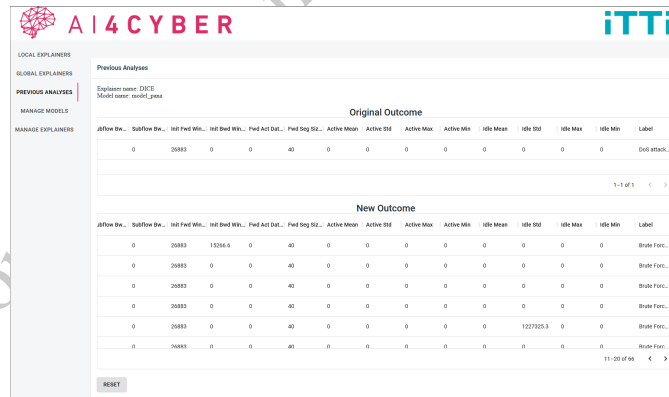


Figure 3. DICE explanation example.

4.3 Local Interpretable Model-Agnostic Explanations (LIME)

LIME (Local Interpretable Model-Agnostic Explanations) is a model-agnostic algorithm that provides explanations based on a generated “local”, simpler model that approximates the behaviour of the original model around a specific input instance. More specifically, in the first step, the instance to be explained is selected, and perturbed versions of the instance are generated by randomly masking or adding noise to the features. Then, the weights of the interpretable features are computed to obtain the importance of each feature in the prediction. This is accomplished by training a linear model on the perturbed instances, where the interpretable features are used as input and the predicted probability of the original model is the output (Ribeiro et al., 2016; Molnar et al., 2020). In the proposed explainer, the explanation is visualized by the plot showing feature importance attributed by LIME, as presented in Fig. 4.

4.4 Shapley Additive Explanations (SHAP)

SHAP (SHapley Additive exPlanations) is a model-agnostic machine learning explanation method based on assigning value points to each input feature to explain the prediction. It is done by showing to which

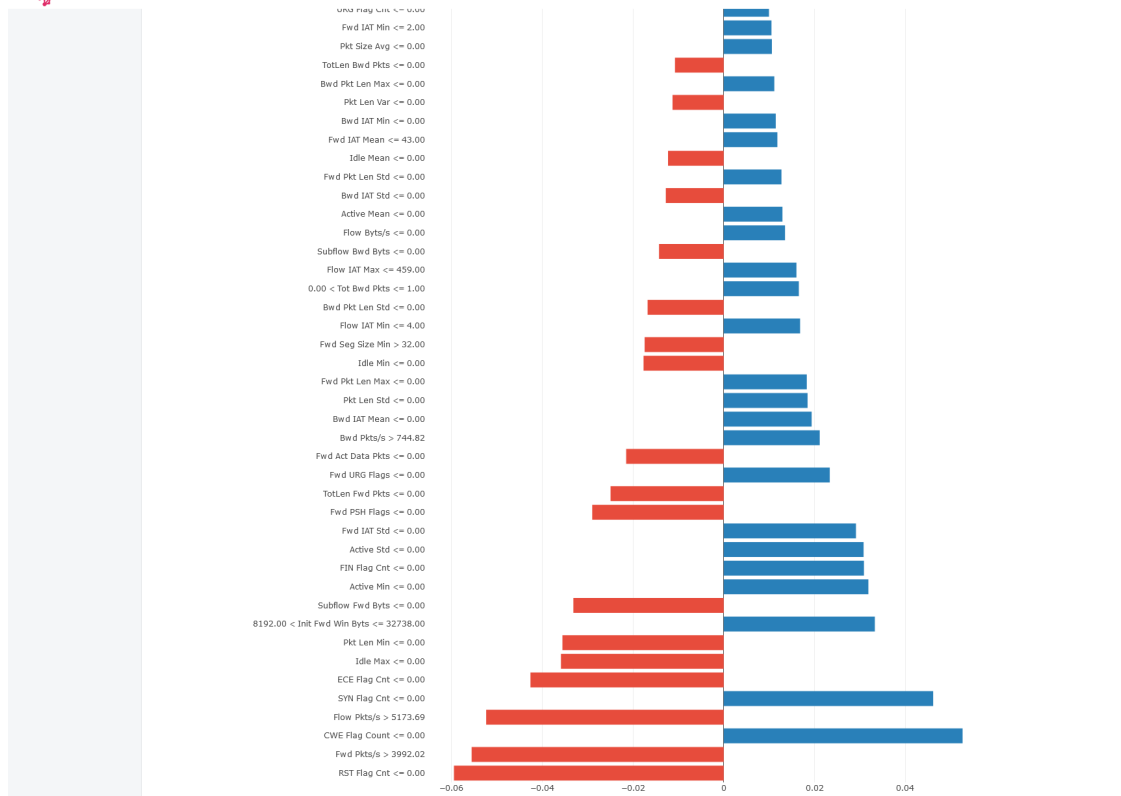


Figure 4. LIME explanation example.

142 extent a specific feature contributes to the final prediction. The SHAP score represents the difference
 143 between expected predictions when the feature exists and when the feature does not exist, calculated on
 144 the average for all the possible combinations of features (Molnar et al., 2020; Lundberg and Lee, 2017).
 145 Complete interpretation of the prediction is visualized by displaying these values for each feature on a bar
 146 or summary plot, as shown in Fig. 5.

147 4.5 Tree-based explainer

148 The decision tree method is a model-agnostic technique to provide model explanations by constructing a
 149 tree-like model of decisions and their potential repercussions. The algorithm ingests the sample along with
 150 the classification performed by the AI model and compares them with pre-trained explanation clusters to
 151 find the closest match. Each node in the tree indicates a choice based on a certain data aspect or attribute,
 152 and the edges reflect the various consequences of that decision (Quinlan, 1986). The decision tree linked to
 153 the closest cluster is used to explain the model prediction (Fig. 6).

154 5 GLOBAL EXPLAINERS

155 5.1 Accumulated Local Effects (ALE)

156 ALE (Accumulated Local Effects) method to visualize and interpret model behaviour has been introduced
 157 in (Apley and Zhu, 2020). This algorithm computes the model predictions for data points within intervals of
 158 a particular feature. After calculating predictions for different intervals, while keeping other variables fixed,
 159 the algorithm calculates differences between adjacent predictions, accumulates the effects and plots them.
 160 The visualized output for a given range of a specific feature allows for understanding the influence of this
 161 feature in predictions (Fig. 7). In particular, this algorithm can be useful to detect non-linear relationships
 162 between variables and the model output.

163 5.2 Partial Dependence Plot (PDP)

164 PDP (Partial Dependence Plot) is a technique to visualize the models predicted outcome, based on the
 165 feature changes. In general, the approach is an alternative to the ALE technique. However, instead of
 166 accumulating local effects, the visualization is provided based on the average of the model predictions over
 167 a dataset consisting of selected features in a specific range (Friedman, 2001). PDP helps in understanding

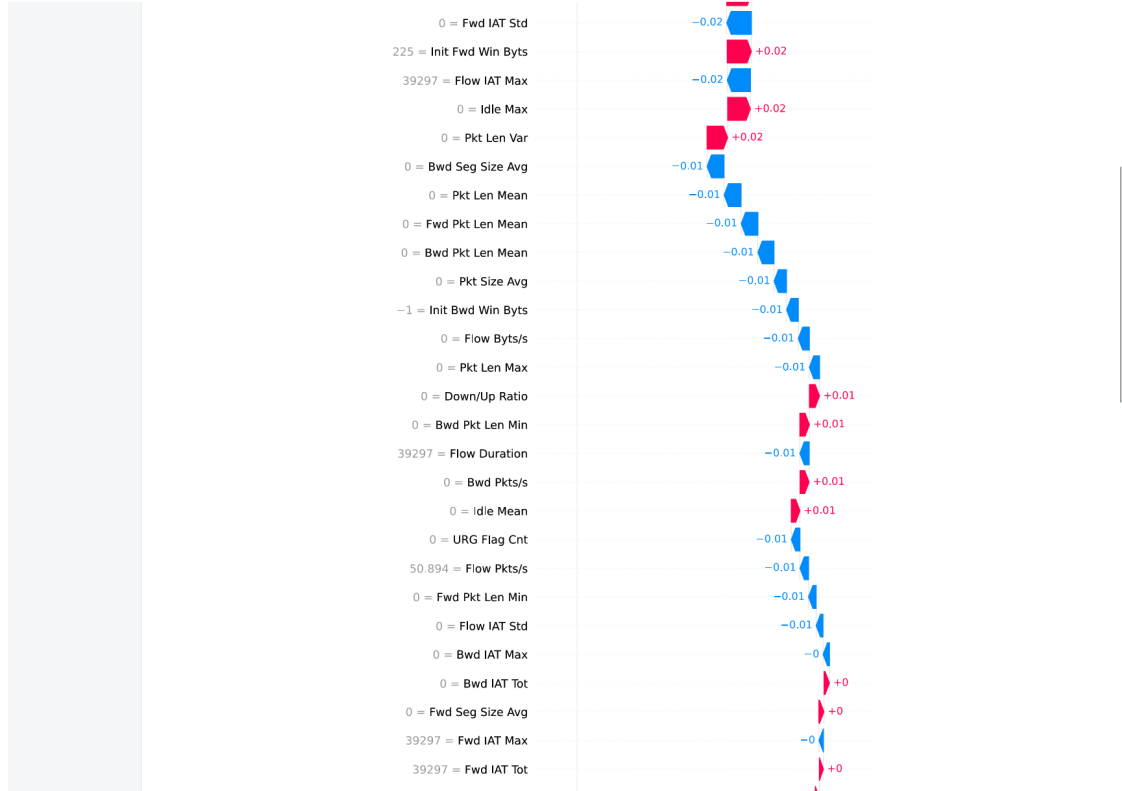


Figure 5. SHAP explanation example.

the global effect of the given feature, ignoring interaction with other features and is useful to interpret complex models by providing a straightforward representation of the feature influence on the predicted output, as presented in Fig. 8.

5.3 Individual Conditional Expectation (ICE)

ICE (Individual Conditional Expectation) plots are another technique used for understanding the relationship between input variables and model predictions, similar to PDP and ALE plots. ICE plots provide a visualization showing the prediction behaviour for each instance individually. This results in plotting one line per instance as shown in Fig. 9, in contrary to one line overall plotted using the PDP method (Goldstein et al., 2013; Molnar et al., 2020).

In general, the PDP provides visualization of the average of the lines calculated and plotted using the ICE method. In the proposed dashboard, there is also an included explainer combining the ICE and PDP plots (Fig. 10, thus providing a more comprehensive view of how selected features influence the model prediction both on average and for varied individual instances. This can be helpful in highlighting of potential variations and anomalies that might be overlooked using only a single method.

5.4 Permutation Feature Importance (PFI)

PFI (Permutation Feature Importance) is a model-agnostic technique helping in model interpretation based on measurements of the model accuracy by shuffling the feature values. In other words, PFI measures the variation of the prediction error applying permutation of the feature values. Based on this approach, a given feature can be evaluated as important to the model's decisions when shuffling its values increases the model error because in this case, the model relied on the feature for the prediction. The model ignores the feature or the feature is less important when changes in its output do not impact the prediction error (Breiman, 2001; Molnar et al., 2020).

5.5 RuleFit

RuleFit shown in Fig. 12, is a method in which a set of rules generated from decision trees is combined with linear models. The final model is a linear combination of these rules and original features making predictions based on both. To identify the non-linear relationships in the data, the RuleFit algorithm first constructs a decision tree ensemble, often a random forest. Using a process known as rule extraction, the

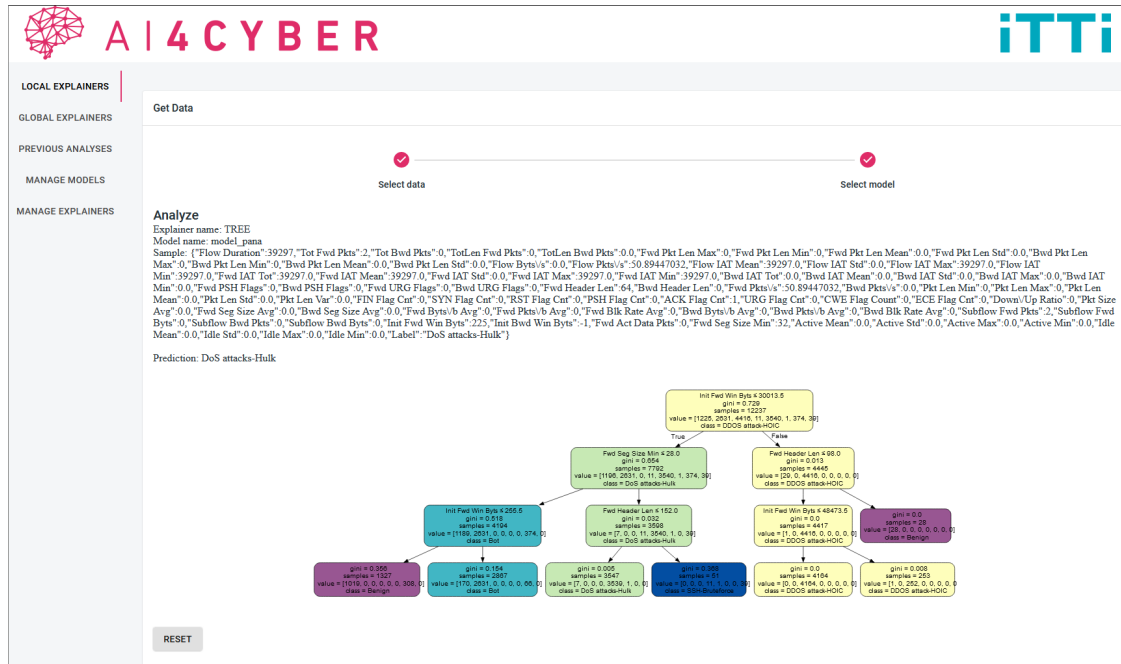


Figure 6. Decision tree explanation example.

195 decision tree algorithm is then turned into a collection of rules. To develop a hybrid model, the extracted
 196 rules are merged with linear models such as linear regression or logistic regression. The linear models are
 197 utilized to represent the linear patterns of the data, while the non-linear correlations are captured by the
 198 retrieved rules from the decision tree ensemble (Friedman and Popescu, 2008; Molnar et al., 2020). The
 199 visualization provides insights not only into the feature importance but also shows which different feature
 200 combinations affect the model output.

201 6 CONCLUSIONS

202 In this paper, a comprehensive explainability module for AI algorithms, called TRUST4AI.XAI has
 203 been proposed. The tool has been developed within the AI4CYBER project to provide AI explainability
 204 capabilities. The paper presents 10 different explanation and visualization techniques, namely: Anchors,
 205 DICE, LIME, SHAP, Decision Tree-based explanations, ALE, PDP, ICE, PFI and RuleFit, implemented as
 206 the backend in the tool presented, and available for the user via the intuitive dashboard. Upon setup of
 207 the tool, these methods are made available for anyone, regardless of their level of expertise and without
 208 necessitating any coding knowledge. This is achieved by setting up a scalable architecture allowing for

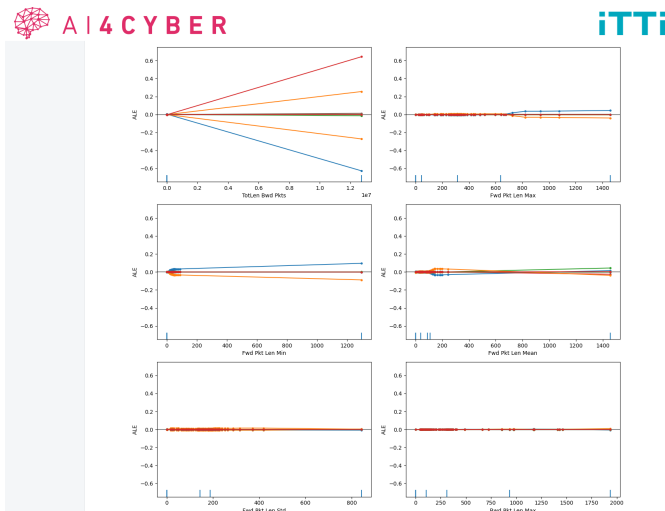


Figure 7. ALE explanation example.

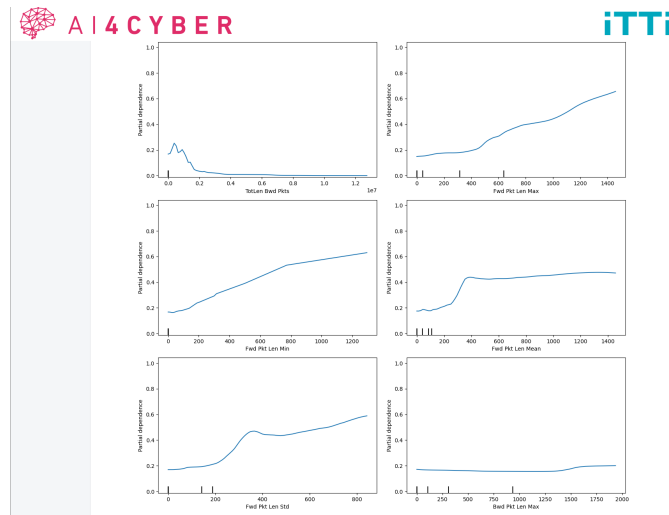


Figure 8. PDP explanation example.

connectivity with various ML models, data preprocessing, and serving the results of the xAI methods to the user in a browser window.

Employing a broad range of different explainability methods on the same model or data sample offers several advantages in comparison to relying only on a single method to interpret the model behaviour. The explanations obtained using different methods provide unique insights into the AI model's decision-making processes, offering a more comprehensive understanding of the underlying logic of the black box.

Parallel use of different explanations can also help to evaluate model robustness, ensuring that the decision-making aligns with domain knowledge. This approach helps to avoid biases that may be observable using one specific method but not another, and facilitates the refining or fine-tuning of the model by allowing a detailed and comprehensive analysis of particular feature importance and its contributions to the model output. The tool is being rolled out to the end users in the AI4Cyber project, and a systematic evaluation of its usability and effectiveness is planned to inform iterative enhancements based on empirical user data.

ACKNOWLEDGMENTS

This research is funded under the Horizon Europe AI4CYBER Project, which has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101070450

REFERENCES

- (2024). Ai4cyber web page at the european commission's cordis portal. <https://cordis.europa.eu/project/id/101070450>. Accessed April 23, 2024.
- (2024). Trustworthy artificial intelligence for cybersecurity reinforcement and system resilience (ai4cyber). <https://ai4cyber.eu/>. Accessed April 23, 2024.
- Abou El Houda, Z., Brik, B., and Khoukhi, L. (2022). "why should i trust your ids?": An explainable deep learning framework for intrusion detection systems in internet of things networks. *IEEE Open Journal of the Communications Society*, 3:1164–1176.
- Apley, D. W. and Zhu, J. (2020). Visualizing the effects of predictor variables in black box supervised learning models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(4):1059–1086.
- Arreche, O., Guntur, T. R., Roberts, J. W., and Abdallah, M. (2024). E-xai: Evaluating black-box explainable ai frameworks for network intrusion detection. *IEEE Access*.
- Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., and Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58:82–115.
- Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1):5–32.
- Choraś, M., Pawlicki, M., Puchalski, D., and Kozik, R. (2020). Machine Learning – The Results Are Not the only Thing that Matters! What About Security, Explainability and Fairness? BT - Computational

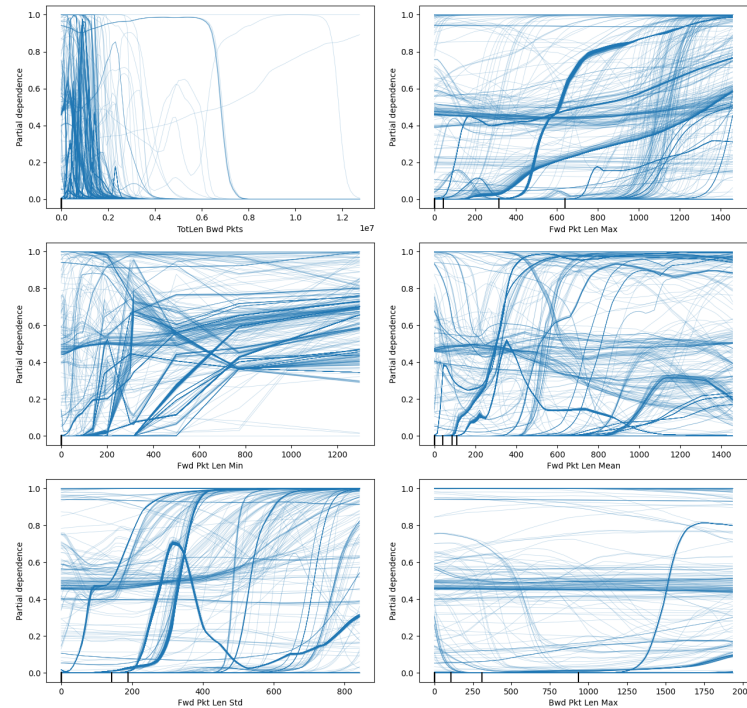


Figure 9. ICE explanation example.

- 245 Science – ICCS 2020. In Krzhizhanovskaya, V. V., Závodszky, G., Lees, M. H., Dongarra, J. J., Sloot, P.
 246 M. A., Brissos, S., and Teixeira, J., editors, , pages 615–628, Cham. Springer International Publishing.
- 247 Damasevicius, R., Venckauskas, A., Grigaliunas, S., Toldinas, J., Morkevicius, N., Aleliunas, T., and
 248 Smuikys, P. (2020). Litnet-2020: An annotated real-world network flow dataset for network intrusion
 249 detection. *Electronics*, 9(5):800.
- 250 Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of*
 251 *Statistics*, 29(5).
- 252 Friedman, J. H. and Popescu, B. E. (2008). Predictive learning via rule ensembles.
- 253 Goldstein, A., Kapelner, A., Bleich, J., and Pitkin, E. (2013). Peeking Inside the Black Box: Visualizing
 254 Statistical Learning with Plots of Individual Conditional Expectation.
- 255 Hariharan, S., Rejimol Robinson, R., Prasad, R. R., Thomas, C., and Balakrishnan, N. (2023). Xai
 256 for intrusion detection system: comparing explanations based on global and local scope. *Journal of*
 257 *Computer Virology and Hacking Techniques*, 19(2):217–239.
- 258 Hu, B., Tunison, P., Vasu, B., Menon, N., Collins, R., and Hoogs, A. (2021). Xaitk: The explainable ai
 259 toolkit. *Applied AI Letters*, 2(4):e40.
- 260 Kaur, R., Gabrijelčič, D., and Klobučar, T. (2023). Artificial intelligence for cybersecurity: Literature
 261 review and future research directions. *Information Fusion*, 97:101804.
- 262 Lundberg, S. M. and Lee, S.-I. (2017). A Unified Approach to Interpreting Model Predictions. In Guyon,
 263 I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors,
 264 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- 265 Mane, S. and Rao, D. (2021). Explaining network intrusion detection system using explainable ai
 266 framework. *arXiv preprint arXiv:2103.07110*.
- 267 Molnar, C., Casalicchio, G., and Bischl, B. (2020). Interpretable machine learning—a brief history, state-of-
 268 the-art and challenges. In *Joint European conference on machine learning and knowledge discovery in*
 269 *databases*, pages 417–431. Springer.
- 270 Mothilal, R. K., Sharma, A., and Tan, C. (2020). Explaining machine learning classifiers through diverse
 271 counterfactual explanations. In *Proceedings of the 2020 conference on fairness, accountability, and*
 272 *transparency*, pages 607–617.
- 273 Neupane, S., Ables, J., Anderson, W., Mittal, S., Rahimi, S., Banicescu, I., and Seale, M. (2022). Explain-
 274 able intrusion detection systems (x-ids): A survey of current methods, challenges, and opportunities.
 275 *IEEE Access*, 10:112392–112415.
- 276 Pawlicka, A., Choraś, M., and Pawlicki, M. (2020). Cyberspace threats: not only hackers and criminals.

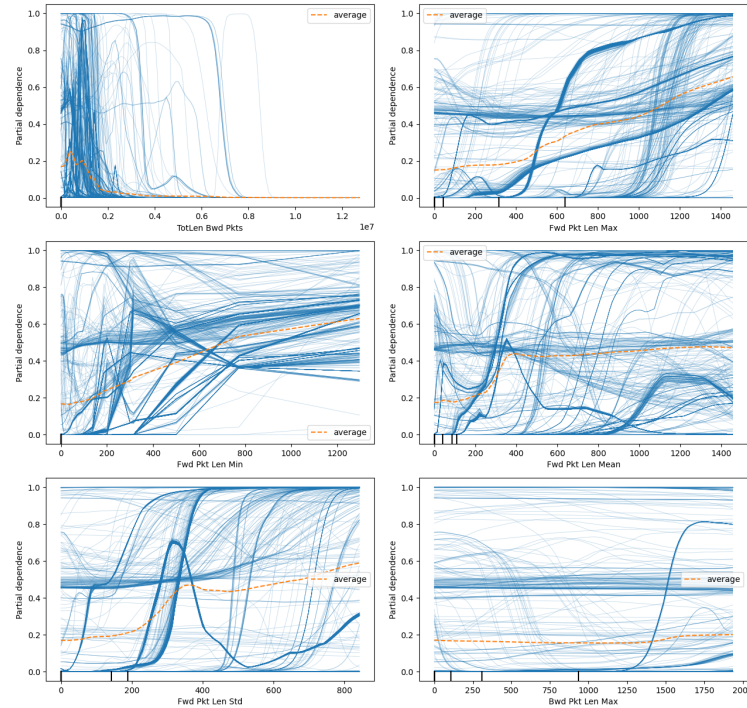


Figure 10. Combination of ICE and PDP explanation – an example.

- raising the awareness of selected unusual cyberspace actors-cybersecurity researchers' perspective. In *Proceedings of the 15th International Conference on Availability, Reliability and Security*, pages 1–11.
- Pawlicka, A., Choraś, M., and Pawlicki, M. (2021). The stray sheep of cyberspace aka the actors who claim they break the law for the greater good. *Personal and Ubiquitous Computing*, 25(5):843–852.
- Pawlicki, M., Pawlicka, A., Kozik, R., and Choraś, M. (2023). The survey and meta-analysis of the attacks, transgressions, countermeasures and security aspects common to the Cloud, Edge and IoT. *Neurocomputing*, 551:126533.
- Pawlicki, M., Pawlicka, A., Kozik, R., and Choraś, M. (2024). Advanced insights through systematic analysis: Mapping future research directions and opportunities for xAI in deep learning and artificial intelligence used in cybersecurity. *Neurocomputing*, page 127759.
- Quinlan, J. R. (1986). Induction of decision trees. *Machine learning*, 1:81–106.
- Rafy, M. F. (2024). *Artificial Intelligence in Cyber Security*.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2018). Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Spinner, T., Schlegel, U., Schäfer, H., and El-Assady, M. (2019). explainer: A visual analytics framework for interactive and explainable machine learning. *IEEE transactions on visualization and computer graphics*, 26(1):1064–1074.
- Warnecke, A., Arp, D., Wressnegger, C., and Rieck, K. (2020). Evaluating explanation methods for deep learning in security. In *2020 IEEE european symposium on security and privacy (EuroS&P)*, pages 158–174. IEEE.
- Xu, Y., Liu, X., Cao, X., Huang, C., Liu, E., Qian, S., Liu, X., Wu, Y., Dong, F., Qiu, C.-W., Qiu, J., Hua, K., Su, W., Wu, J., Xu, H., Han, Y., Fu, C., Yin, Z., Liu, M., Roepman, R., Dietmann, S., Virta, M., Kengara, F., Zhang, Z., Zhang, L., Zhao, T., Dai, J., Yang, J., Lan, L., Luo, M., Liu, Z., An, T., Zhang, B., He, X., Cong, S., Liu, X., Zhang, W., Lewis, J. P., Tiedje, J. M., Wang, Q., An, Z., Wang, F., Zhang, L., Huang, T., Lu, C., Cai, Z., Wang, F., and Zhang, J. (2021). Artificial intelligence: A powerful paradigm for scientific research. *The Innovation*, 2(4):100179.
- Yang, W., Le, H., Laud, T., Savarese, S., and Hoi, S. C. (2022). Omnixai: A library for explainable ai. *arXiv preprint arXiv:2206.01612*.

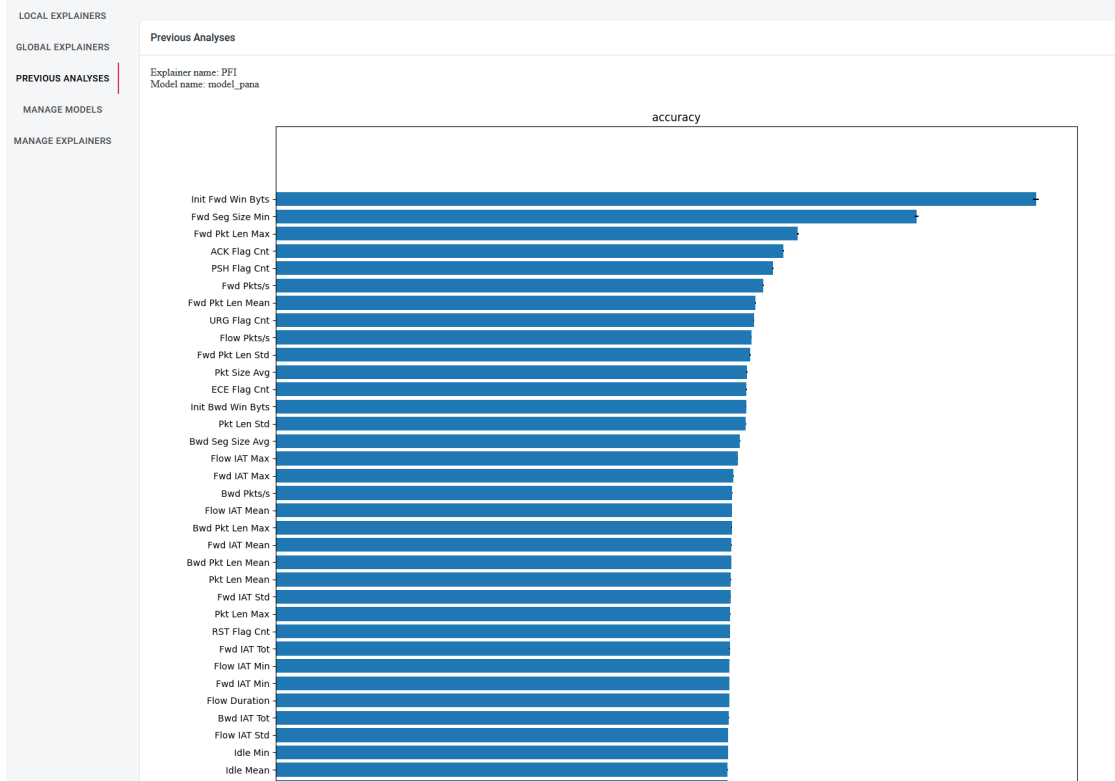


Figure 11. PFI explanation example.

LOCAL EXPLAINERS

GLOBAL EXPLAINERS

PREVIOUS ANALYSES

MANAGE MODELS

MANAGE EXPLAINERS

Previous Analyses

Explainer name: RULEFIT
Model name: model_pana

rule	type	coef	support	importance
PSH Flag Cnt	linear	-0.4040504517644865	1	0.197634953236349...
Idle Min <= 58030570.0 & Bwd IAT Max <= 10007699.0 & Idle Mean <= 57999088.0	rule	1.0359088898215307	0.9640151515151515	0.1929404686854419
Bwd IAT Max <= 953133.5 & TotLen Fwd Pkts <= 1686.0 & Flow Duration <= 1792613.5	rule	-0.5374432148678372	0.8522727272727273	0.190700670517184...
Fwd IAT Mean > 194.0 & Init Fwd Win Byts <= 15425.5 & Fwd Pkt Len Std <= 208.1890106201172 & Fwd Seg Size Min <= 22.0	rule	-0.46050286886332...	0.2108585858585858...	0.187847492746323...
Bwd IAT Max <= 953169.5 & Fwd Pkts/s <= 30.82939052581787	rule	-0.5729436907168711	0.0959595959595959...	0.168752442326640...
Fwd Pkts/s <= 45983.087890625 & Init Fwd Win Byts <= 20741.5 & Fwd Seg Size Min > 22.0	rule	0.443869678118984...	0.1666666666666666...	0.165420462237499...
Subflow Bwd Byts > 1912.0 & ECE Flag Cnt > 0.5	rule	1.6107146423098666	0.0101010101010101...	0.161063246928008...

21 - 40 of 233

RESET

Figure 12. RuleFit explanation example.