

International Journal of Advanced Research in Science, Communication and Technology (IJARSCT)

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 2, March 2025

Interpretable AI: Enhancing Transparency and Fairness in Decision-Making

Rajeev Nair¹, Rosemol Thomas¹, Sandra MV¹, Ms. Siji K B²

Students, MCA, Vidya Academy of Science and Technology Thalakkottukara, Thrissur, India¹ Assistant Professor, Department of Computer Applications² Vidya Academy of Science and Technology, Thalakkottukara, Thrissur, India

Abstract: Interpretable Artificial Intelligence aims to make machine learning models more transparent, interpretable, and accountable, addressing the "black box" nature of traditional AI systems. As AI plays a critical role in high-stakes domains like healthcare, finance, and autonomous systems, ensuring trust and fairness in decision-making has become essential and this paper also explores key techniques in AI. This study adopts a mixed-methods approach to analyse, evaluate, and compare XAI techniques across key domains. This research examines a three-phase approach in XAI, focusing on exploring different methods, evaluating their impact in real- world applications, and analysing the trade-offs between interpretability and model performance. XAI enhances transparency, bias detection, and user trust, but still it faces challenges, such as the trade-off between interpretability and accuracy, as well as computational complexity. Future studies focus on improving model interpretability, enhancing human-AI interaction, and promoting fairness in AI-driven decisions.ncy.

Keywords: Explainable Artificial Intelligence , Transparency & Trust , LIME & SHAP , Model-Agnostic Methods, Gradient-Based Methods , Propagation-Based Methods , Meta- Explanations , Accuracy vs. Simplicity , High-Stakes Domains

I. INTRODUCTION

Explainable Artificial Intelligence (XAI) is a transformational approach for the AI systems mainly aims at growing needs of the transparency and accountability. Due to increase in AI systems usage, more complex AI models are not understandable for the human beings. Therefore, XAI helps to bridge the gap between complex AI models and human comprehension by providing clear, interpretable explanations for predictions and decisions. Hence, XAI plays a crucial role in making AI more transparent, responsible, and aligned with human values. XAI represents a crucial development in the field of artificial intelligence, aiming to make complex AI models more transparent, trustworthy, and accountable. XAI addresses this limitation by offering human-understandable explanations for AI-generated outcomes, empowering users to trust, evaluate, and improve AI systems. XAI emerges as a solution to the long-standing challenge of " black box" AI models, which produce decisions without revealing how they were made. XAI not only builds user trust but also ensures ethical AI deployment, aids in bias detection, and supports regulatory compliance, helping the way for a future where AI systems are both powerful and responsible.

II. WHY EXPLAINABLE AI?

Explainable AI (XAI) is essential for fostering trust, transparency, and accountability in AI-driven systems. As AI becomes a cornerstone of decision-making in critical areas such as healthcare, finance, security, and autonomous systems, understanding how these systems arrive at their conclusions is vital. Traditional AI models, often described as "black boxes," produce highly accurate results but lack transparency. This may lead to issues of bias, discrimination, or critical errors going undetected. XAI addresses these challenges by providing human-understandable insights into how AI systems process information, identify patterns, and arrive at conclusions. Transparency and trust are crucial for the effectiveness and acceptance of AI, ensuring responsible deployment.

Copyright to IJARSCT www.ijarsct.co.in DOI: 10.48175/IJARSCT-23772







Fig 1. Explainable AI

III. TRUST AND TRANSPARENCY THROUGH XAI

Trust and Transparency are two fundamental components of XAI that determine the effectiveness and acceptance of AI systems, especially in high-stakes domains like healthcare, finance, and autonomous vehicles.

Trust is a user-confidence which helps the user or human- being to choose an AI decision for their particular system. This user-confidence enables the user to choose a proper AI decision for the AI system. Consistency, Accuracy of explanation, simplicity, complexity etc. are some of the key factors of XAI Trust.

Transparency in XAI refers to the extent to which users can understand and trace the decision-making process of a machine learning model. It is crucial for ensuring that users know how and why a model arrives at its conclusions. Model transparency, Explainability of features & Ethical and Fair-Decision making are some of the key concepts or features of XAI transparency.

Trust and transparency are closely inter-related in XAI. Transparency in how models work and make decisions is a necessary foundation for trust. If a user cannot understand how an AI system arrived at a particular decision, they are unlikely to trust that decision, no matter how accurate it is. On the other hand, simply providing transparent explanations without ensuring that they are accurate, consistent, and understandable might lead to confusion and mistrust.

To wrap things up, trust and transparency are essential for making AI systems more understandable, ethical, and reliable. The goal is to ensure that users feel confident not only in the system's predictions but also in the rationale behind them, fostering trust and enabling responsible AI deployment.

IV. METHODS OF XAI

Explainable AI (XAI) methods can be broadly classified into Gradient-Based Methods, Model-Agnostic Methods, Propagation-Based Methods & Meta-Explanations.

1. Model-Agnostic Method

Model-agnostic methods are flexible techniques that can be applied to any type of AI model, regardless of its complexity or structure. These methods provide post-hoc explanations by analyzing the model's input-output relationship without requiring access to the internal workings. It includes techniques like LIME (Local Interpretable Model-agnostic Explanations) & SHAP (Shaply Additive Explanations).

Copyright to IJARSCT www.ijarsct.co.in DOI: 10.48175/IJARSCT-23772



IJARSCT



International Journal of Advanced Research in Science, Communication and Technology (IJARSCT)

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 2, March 2025

2. Gradient-Based Method

Gradient-Based Methods rely on gradients (partial derivatives) to explain the relationship between input features and the output of the AI model. Gradients indicate how small changes in input features affect the prediction. It includes of GRAD- CAM and Integrated Gradients.

3. Propagation-Based Method

Propagation-based methods focus on how information flows through the layers of neural networks to understand feature importance or decision paths. It includes methods like Layer-wise Relevance Propagation & Deep LIFT(Deep Learning Important Features).

4. Meta-Explanations

Meta-explanations go beyond direct feature importance or gradient analysis by offering high-level, generalized insights about the behavior and fairness of AI models. They evaluate how explanations themselves are generated, aggregated, or applied. It includes Explanation Aggregation, Fairness and Bias Analysis & Concept-Based Explanations These approaches collectively enhance the interpretability and transparency of AI models, ensuring that both technical users and non-experts can better trust and validate AI decisions.

V. KEY TECHNIQUES IN XAI: LIME & SHAP

SHAP (SHapley Additive exPlanations) is a technique based on game theory, specifically Shapley values, which determine the contribution of each feature to the model's prediction. SHAP provides both global and local explanations, offering insights into overall feature importance as well as how individual features affect specific predictions. The key advantage of SHAP is its mathematical consistency—each feature's contribution is additive, meaning the sum of the SHAP values equals the model's output. This makes SHAP particularly suitable for tree-based models like XGBoost and LightGBM, though it can be computationally expensive when working with large datasets or models with many features. Despite this, SHAP's rigorous approach is valuable for understanding the full decision-making process of a model in a clear, consistent manner.

LIME (Local Interpretable Model-agnostic Explanations), on the other hand, focuses on providing local explanations for individual predictions by approximating the complex model with an interpretable surrogate model around a specific instance. LIME generates perturbed samples of the instance in question and trains a simpler model, like a linear regression or decision tree, on these samples to approximate the black- box model's behavior in that local region. This technique is model-agnostic, meaning it can be applied to any machine learning model, and it is particularly valued for its simplicity and ease of understanding. However, LIME only provides explanations for individual instances, not the entire model, and its accuracy depends on how well the surrogate model approximates the original model in the local area. While LIME is useful for understanding specific predictions, it lacks the global perspective that SHAP offers.

VI. APPLICATIONS OF XAI

The applications of Explainable AI (XAI) span across various industries and fields where transparency, interpretability, and accountability are critical. Some of the key factors are:

Healthcare is one of the major aspects of Explainable AI (XAI). Healthcare plays a crucial role in improving diagnostic accuracy and supporting clinical decision-making. AI models are used to analyze medical data such as imaging, lab results, and patient history to assist healthcare professionals in diagnosing diseases, predicting outcomes, and recommending treatments. XAI ensures transparency by providing understandable explanations of how these AI-driven decisions are made, which helps doctors trust the system's recommendations. For example, a model used to diagnose cancer from medical images can help to influence its prediction, allowing clinicians to validate the results.

Finance in XAI enhances trust and accountability in decision-making processes such as credit scoring, risk assessment, and fraud detection. Mostly, AI systems are commonly used by banks and lenders to evaluate the credits of individuals and businesses based on historical data, income levels, and spending behaviors. With XAI customers and regulators

Copyright to IJARSCT www.ijarsct.co.in DOI: 10.48175/IJARSCT-23772



440



International Journal of Advanced Research in Science, Communication and Technology (IJARSCT)

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 2, March 2025

can understand why certain decisions are made such as why a loan application was denied or why a transaction was flagged as fraudulent.

Autonomous Driving in XAI are essential for ensuring the safety, accountability, and public acceptance of self-driving cars. Autonomous vehicles rely on AI systems to make real- time decisions based on sensor data, traffic patterns, and environmental conditions which helps to avoid accidents and any cases of emergency situations to avoid damage.

VII. CASE STUDY: HEALTHCARE

Artificial Intelligence (AI) is increasingly being used in healthcare for disease diagnosis, treatment recommendations, and patient care. However, AI models, particularly deep learning, are often seen as "black boxes," making it difficult to understand their decision-making process. Explainable AI (XAI) aims to provide transparency, trust, and interpretability in AI-driven healthcare applications.

Diabetes is a chronic medical condition that affects millions of people worldwide. Early diagnosis and personalized treatment plans can help prevent severe complications like heart disease and kidney failure. Machine learning models are being used to predict the risk of diabetes based on various patient data, such as age, gender, family history, lifestyle, and medical records.

However, doctors and healthcare providers need to trust these predictions, especially when they're making decisions on treatments or lifestyle changes for patients. A machine learning model that simply predicts "high risk" or "low risk" without providing any context or reasoning may not be useful in real-world medical settings.

In a hospital, healthcare professionals can use the AI model to help assess patient risk for diabetes. Once a prediction is made, the system generates an explanation. For instance, for a patient named John, who is overweight, sedentary, and has a family history of diabetes, the model might predict a high risk of developing diabetes in the next 5 years.

VIII. THREE-PHASE APPROACH TO XAI EVALUVATION

This study uses a mixed-methods approach, which means it combines both qualitative (subjective, like opinions) and quantitative (objective, like data and numbers) methods. The goal is to understand how XAI can make machine learning models more transparent and trustworthy. The methodology is divided into three main phases:

Theoretical Analysis of XAI Techniques mainly focuses on two main types: model-agnostic and model-specific methods. Model-agnostic techniques, like LIME and SHAP, can be applied to any machine learning model, offering general explanations for predictions. In contrast, model-specific techniques, such as Layer-wise Relevance Propagation (LRP), are designed for particular models, especially deep learning models, to provide insights into how certain input features influence the model's decisions. Both types aim to enhance the interpretability and trustworthiness of AI systems.

In **Empirical Evaluation** of XAI Methods, the study ap- plies various XAI methods to real-world applications across three distinct domains: healthcare, finance, and autonomous systems. These domains were chosen due to their complexity and the critical need for transparency in AI-driven decision- making. For each domain, the study evaluates the performance of the XAI techniques using several key metrics. These include model accuracy, user satisfaction and trust in AI. By collecting these metrics, the study aims to understand the strengths and weaknesses of different XAI methods in providing transparent, reliable, and actionable insights in healthcare, finance, and autonomous systems.

In **Comparative Analysis**, the study will compare the XAI methods used across the three domains—healthcare, finance, and autonomous systems—focusing on three key factors. Explanation Quality will evaluate how clear, relevant, and understandable the explanations are, based on expert feedback and user studies. Impact on Trust will assess how the explanations influence users' trust in the AI model, measured through surveys that capture users' confidence before and after receiving the explanations. Lastly, Model Performance will analyze the trade-off between the model's interpretability and its predictive accuracy, comparing how well the model performs after applying XAI methods. This comprehensive comparison will provide insights into the effectiveness of different XAI techniques across diverse applications.

Copyright to IJARSCT www.ijarsct.co.in DOI: 10.48175/IJARSCT-23772



441

IJARSCT



International Journal of Advanced Research in Science, Communication and Technology (IJARSCT)

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 2, March 2025

IX. CHALLENGES

Explainable AI (XAI) faces a number of significant challenges that increase in widespread adoption and effectiveness. One of the main challenges is the trade-off between model complexity and interpretability. Complex models, such as deep neural networks, often achieve higher accuracy but are difficult to interpret. In order to make these models more understandable can lead to a reduction in their predictive performance, creating a balance between transparency and accuracy. Additionally, there is no universally accepted metric to evaluate the quality of explanations, making it difficult to access the effectiveness of XAI methods. The subjectivity of human understanding adds to the complexity, as different users, depending on their expertise, may interpret explanations differently, which complicates the design of universally clear explanations.

Accuracy vs Simplicity one challenge in XAI is balancing accuracy and interpretability. Complex AI models, like deep learning, work well and give accurate results, but they're hard to explain. Sometimes, making these models simpler to understand can lower their performance, creating a problem for users who need both accuracy and clear explanations. Different understanding of user another issue is that different people might understand explanations in different ways. What makes sense to an expert may not be clear to someone without technical knowledge. So the researcher works to make the explanation which suit for the user level of understanding.

In XAI, ensuring that explanations are stable across similar inputs is crucial. In some cases, the explanations provided by XAI techniques can be highly sensitive to small changes in input data, leading to inconsistent or unstable explanations. For example, a slight perturbation in the input could cause radically different explanations, which may confuse users and reduce trust in the system.

X. CONCLUSION

Explainable Artificial Intelligence (XAI) is essential for making AI systems more transparent, ethical, and trustworthy, especially in high-impact domains like healthcare, finance, and autonomous systems. By providing humanunderstandable explanations, XAI enhances user confidence, supports regulatory compliance, and mitigates risks associated with bias and discrimination. While challenges such as balancing accuracy with interpretability and ensuring stable explanations persist, continuous advancements in XAI methodologies are improving AI transparency. Ultimately, XAI plays a crucial role in shaping the future of responsible AI, ensuring that AI- driven decisions are not only accurate but also fair, reliable, and aligned with human values.

REFERENCES

- [1]. Alvarez-Melis, D., & Jaakkola, T. S. (2017). A causal framework for explaining the predictions of black-box sequence-to-sequence models. CoRR, abs/1707.01943.
- [2]. Andrews, R., Diederich, J., & Tickle, A. B. (1995). Survey and critique of techniques for extracting rules from trained artificial neural networks. Knowledge-Based Systems, 8(6), 373–389.
- [3]. Angelov, P. P., & Gu, X. (2018). Toward anthropomorphic machine learning. Computer, 51, 18–27.
- [4]. Arrieta, A. B., D'ıaz-Rodr'ıguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garc'ıa, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward re- sponsible AI. Information Fusion, 58, 82–115.
- [5]. Bach, S., Binder, A., Montavon, G., Klauschen, F., Mu⁻Iler, K.R., Samek, W.: On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS ONE 10(7), e0130140 (2015).
- [6]. Bau, D., Zhou, B., Khosla, A., Oliva, A., Torralba, A.: Network dissection: Quantifying interpretability of deep visual representations. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- **[7].** pp. 6541–6549 (2017).
- [8]. Baehrens, D., Schroeter, T., Harmeling, S., Kawanabe, M., Hansen, K., Mu[°]ller, K.R.: How to explain individual classification decisions. Journal of Machine Learning Research 11, 1803–1831 (2010).
- [9]. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: International Conference on Learning Representations (ICLR). (2015).

Copyright to IJARSCT www.ijarsct.co.in DOI: 10.48175/IJARSCT-23772



IJARSCT



International Journal of Advanced Research in Science, Communication and Technology (IJARSCT)

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 2, March 2025

- [10]. Cires, an, D., Meier, U., Masci, J., Schmidhuber, J.: A committee of neural networks for traffic sign classification. In: International Joint Conference on Neural Networks (IJCNN). pp. 1918–1921 (2011).
- [11]. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 248–255 (2009).
- [12]. Doshi-Velez, F., Kim, B.: Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608 (2017).
- [13]. Doyle, D., Tsymbal, A., & Cunningham, P. (2003). A review of explanation and explanation in case-based reasoning (Technical Report). Dublin: Trinity College Dublin, Department of Computer Science.
- [14]. Dressel, J., & Farid, H. (2018). The accuracy, fairness, and limits of predicting recidivism. Science Advances, 4, eaao5580.
- [15]. Eitel, F., Soehler, E., Bellmann-Strobl, J., Brandt, A.U., Ruprecht, K., Giess, R.M., Kuchling, J., Asseyer, S., Weygandt, M., Haynes, J.D., et al.: Uncovering convolutional neural network decisions for diagnosing multiple sclerosis on conventional mri using layer-wise relevance prop- agation. arXiv preprint arXiv:1904.08771 (2019).
- [16]. Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The PASCAL Visual Object Classes Challenge: A Retrospective. International Journal of Computer Vision 111(1), 98–136 (2015).
- [17]. Fabian, B., Ermakova, T., & amp; Junghanns, P. (2015). Collaborative and secure sharing of healthcare data in multi-clouds. Information Systems, 48, 132-150.
- [18]. Felzmann H, Villaronga EF, Lutz C, Tamo'Larrieux A. Transparency you can trust: Transparency requirements for artificial intelligence between legal norms and contextual concerns. Big Data and Society. 2019;6(1):1–14.
- [19]. Kambatla, K., Kollias, G., Kumar, V., & Grama, A. (2014). Trends in big data analytics. Journal of parallel and distributed computing, 74(7), 2561-2573.
- [20]. Konda SR. Ensuring Trust and Security in AI: Challenges and Solutions for Safe Integration. International Journal of Computer Science and Technology. 2019;3(2):71–86.
- [21]. Jones, S. L., and Shah, P. P. (2016). Diagnosing the locus of trust: a temporal perspective for trustor, trustee, and dyadic influences on perceived trustworthiness. J. Appl. Psychol. 101, 392–414. doi: 10.1037/apl0000041.
- [22]. Kaplan, A. D., Kessler, T. T., Brill, J. C., and Hancock, P. A. (2021). Trust in artificial intelligence: Metaanalytic findings. Hum. Factors 65, 337–359. doi: 10.1177/00187208211013988
- [23]. Kulms, P., and Kopp, S. (2018). A social cognition perspective on human-computer trust: the effect of perceived warmth and competence on trust in decision-making with computers. Front. Digit. Humanit. 5:14. doi: 10.3389/fdigh.2018.00014
- [24]. Lee, M. K., and Rich, K. (2021). Who is included in human perceptions of AI?: trust and perceived fairness around healthcare AI and cultural mistrust. In: Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, 1–14.
- [25]. Pavithra, B., Niranjanamurthy, M., Smitha, G. V., Kiran, R., & Chan- drika, M. (2023, May). Transactional and Sequential model for pre- processing the textual data. In 2023 4th International conference on intelligent engineering and management (ICIEM) (pp. 1-6). IEEE
- [26]. Muir, B. M. (1987). Trust between humans and machines, and the design of decision aids. Int. J. Man Mach. Stud. 27, 527–539. doi: 10.1016/S0020-7373(87)80013-5.

