

Review Article

Open Access

Bridging the Gap: Enhancing Trust and Transparency in Machine Learning with Explainable AI

Pushkar Mehendale

San Francisco, CA, USA

ABSTRACT

Explainable Artificial Intelligence (XAI) aims to address the complexity and opacity of AI systems, often referred to as "black boxes." It seeks to provide transparency and build trust in AI, particularly in domains where decisions impact safety, security, and ethical considerations. XAI approaches fall into three categories: opaque systems that offer no explanation for their predictions, interpretable systems that provide some level of justification, and comprehensible systems that enable users to reason about and interact with the AI system. Automated reasoning plays a crucial role in achieving truly explainable AI. This paper presents current methodologies, challenges, and the importance of integrating automated reasoning for XAI. It is grounded in a thorough literature review and case studies, providing insights into the practical applications and future directions for XAI.

*Corresponding author

Pushkar Mehendale, San Francisco, CA, USA.

Received: December 07, 2022; **Accepted:** December 14, 2022; **Published:** December 24, 2022

Keywords: Machine Learning, Explainable AI, Transparency, Interpretability, Artificial Intelligence

Introduction

Artificial Intelligence (AI) systems based on deep learning and advanced machine learning techniques have made remarkable strides, achieving impressive results in various fields. However, these systems often lack transparency, making it challenging for users to comprehend their decision-making processes. This opacity can lead to distrust and skepticism, particularly in critical applications such as healthcare, finance, and autonomous systems. Explainable AI (XAI) emerges as a promising solution to address this issue, aiming to enhance the understandability of AI systems for humans [1].

The fundamental objective of XAI is to develop models that can offer clear and concise explanations for their actions and decisions. To accomplish this, XAI systems can be broadly categorized into three main types: opaque systems, interpretable systems, and comprehensible systems. Opaque systems, at one end of the spectrum, provide no insights into their internal mechanisms, leaving users in the dark about their decision-making processes [2]. Interpretable systems, on the other hand, allow for mathematical analysis of their operations, enabling a deeper understanding of their behavior. Comprehensible systems, occupying a middle ground, leverage emitted symbols or visualizations to provide user-driven explanations, bridging the gap between opaque and interpretable systems [3].

Beyond these three types of XAI systems, we introduce the concept of truly explainable systems. Truly explainable systems make a significant leap forward by integrating automated reasoning to generate human-understandable explanations without requiring extensive human interpretation. This type of XAI system is still in its nascent stages of development, but it holds immense potential to

revolutionize our interactions with AI systems. Truly explainable systems can enhance trust, foster collaboration, and ultimately pave the way for more responsible and ethical AI applications [4, 5].

Literature Review

The demand for Explainable Artificial Intelligence (XAI) arises in domains where comprehending the AI decision-making process is vital. Healthcare is a prime example, as AI systems play a crucial role in diagnosing diseases and proposing treatments. Without explainability, medical professionals may hesitate to rely on these AI recommendations. Similarly, in finance, AI models employed for credit scoring or fraud detection require transparency to guarantee fairness and adherence to regulations [6].

Explainable AI Categories

1. **Opaque Systems:** The lack of transparency and interpretability in certain AI systems, such as proprietary AI systems and some deep learning models, poses a significant challenge. These systems operate as "black boxes," where the input data is transformed into output without any explanation or insight into the decision-making process. This opacity makes it difficult to understand how the system arrives at its conclusions, evaluate its fairness and bias, and ensure that it is functioning as intended. This lack of transparency can hinder trust in AI systems and limit their widespread adoption and acceptance.
2. **Interpretable Systems:** Machine learning models that prioritize interpretability provide users with insight into the relationships between input and output variables. These models, such as decision trees and linear regression, enable users to comprehend how predictions are made through mathematical or logical analysis. However, this focus on interpretability can come at the cost of accuracy, making these models less effective in handling intricate tasks [7].

3. Comprehensible Systems: These systems, which include techniques such as feature importance rankings, saliency maps, and attention mechanisms in neural networks, offer more intuitive explanations by emitting symbols or visualizations that help users understand how a conclusion is reached. However, while they provide a deeper insight into the decision-making process, they don't eliminate the need for user interpretation. Users still need to actively engage with these visualizations to fully comprehend the factors influencing the system's conclusions [8].

Key Concepts in Explainable AI

- 1. Explanation Triggers:** Explanations are crucial in situations where expectations are violated, as they help users comprehend why certain outcomes or events occurred. By identifying the factors that prompt the need for explanation, such as surprising or unexpected results, designers of Explainable Artificial Intelligence (XAI) systems can create more effective and user-centric interfaces. These systems should provide clear and concise explanations that address the user's questions and concerns, thereby enhancing trust and transparency in the interactions between humans and AI systems.
- 2. Self-Explanation:** Encouraging users to generate their explanations for AI systems' predictions can significantly enhance their comprehension and confidence in the systems. By actively engaging in the explanation process, users gain a deeper understanding of how the AI arrives at its conclusions, fostering a sense of transparency and accountability. This promotes trust in the AI's capabilities, allowing users to make more informed decisions and interact with the technology more effectively [9].
- 3. Boundary Conditions:** Understanding the limitations and boundary conditions of an AI system is crucial for fostering appropriate trust and reliance from users. By clearly defining the capabilities and constraints of the AI system, users can develop realistic expectations and make informed decisions about when and how to engage with the system. This transparency enables users to trust the AI system within its specified boundaries and avoid over-reliance or misuse. It also empowers them to recognize situations where human intervention or alternative approaches may be necessary, promoting responsible and effective interactions between humans and AI systems.
- 4. Contrastive Explanations:** Providing explanations that contrast different decisions or outcomes can help users understand the rationale behind a particular choice. By comparing and contrasting alternatives, users can gain insights into the factors that influenced the decision-making process and the potential consequences of each option. This approach enhances comprehension, allowing users to better evaluate the decision and its implications, and make informed choices in similar situations. Additionally, contrasting explanations can uncover hidden assumptions or biases, promoting critical thinking and fostering a deeper understanding of the decision-making process [10, 11].

Methodologies in Explainable AI

Model-Agnostic Approaches

These approaches do not depend on the underlying model and can be applied to any AI system. Examples include LIME (Local Interpretable Model-agnostic Explanations) and SHAP (SHapley Additive exPlanations), which provide local explanations for individual predictions by approximating the model with a simpler, interpretable model around the prediction of interest [3, 4].

1. Local Interpretable Model-agnostic Explanations (LIME): LIME works by perturbing the input data and observing the changes in the output. It creates a new, interpretable model that approximates the predictions of the original model locally around the prediction of interest. This helps in understanding the contribution of individual features to the prediction.

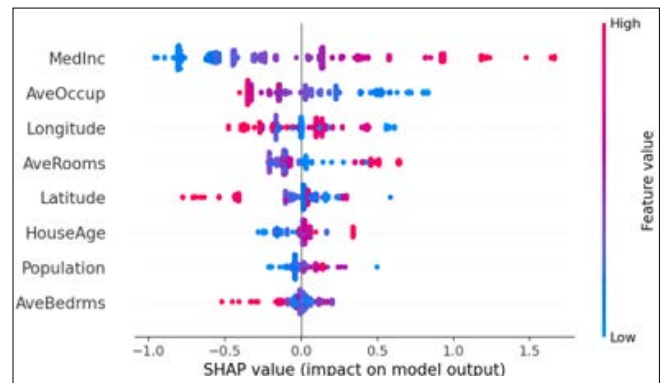


Figure 1: Example of SHAP Value Explanation [12].

2. SHapley Additive exPlanations (SHAP): SHAP values are derived from cooperative game theory and provide a way to fairly distribute the prediction among the features. By calculating the average marginal contribution of each feature to the prediction, SHAP provides a comprehensive explanation of the model's output [13].

Model-Specific Approaches

These approaches are tailored to specific types of models. For instance, decision trees and rule-based systems are inherently interpretable, while neural networks may use techniques like attention mechanisms, layer-wise relevance propagation, or gradient-based methods to provide explanations.

- 1. Attention Mechanisms:** Attention mechanisms in neural networks allow the model to focus on specific parts of the input when deciding. This can be visualized to show which parts of the input were most influential in the decision-making process.
- 2. Layer-Wise Relevance Propagation (LRP):** LRP is a technique used to decompose the prediction of a neural network and attribute the prediction to the input features. It helps in understanding the contribution of each input feature to the final prediction.

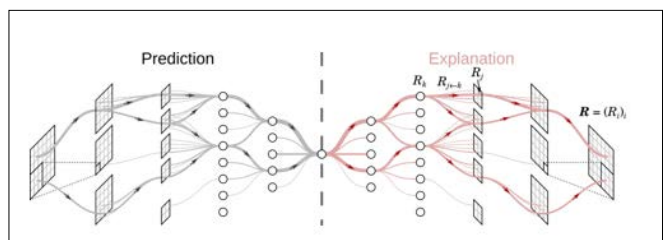


Figure 2: LRP Visualization Example [14].

Visual and Interactive Explanations

Visual explanations, such as saliency maps and heatmaps, highlight the parts of the input that most influenced the AI's decision. Interactive tools allow users to query the AI system and receive explanations in real-time, fostering a more dynamic understanding of the model's behaviour.

- 1. Saliency Maps:** Saliency maps highlight the regions of an input image that are most relevant to the prediction. This

is particularly useful in image classification tasks, where understanding which parts of the image contributed to the prediction can provide insights into the model's behaviour.

- 2. Interactive Tools:** Interactive tools such as IBM's AI Explainability 360 or Microsoft's InterpretML allow users to interact with AI models, query them, and receive explanations. These tools provide a user-friendly interface for exploring the model's decision-making process and understanding its behavior in different scenarios.

Challenges in Explainable AI

Complexity vs. Interpretability

There is often a trade-off between the complexity of the model and its interpretability. More complex models tend to be more accurate but harder to explain. Balancing these two aspects is a significant challenge in the development of XAI systems.

The interplay between the complexity of a model and its interpretability presents a fundamental challenge in the field of Explainable Artificial Intelligence (XAI) systems. The tension between accuracy and explainability is a critical issue that must be addressed to unlock the full potential of XAI [1].

On the one hand, more complex models often achieve higher levels of accuracy. This is because they can capture more nuanced patterns and relationships in the data. However, as models become more complex, they also become more difficult to understand and interpret. This can make it challenging for humans to understand how the model makes its decisions and to trust its predictions.

On the other hand, simpler models are generally easier to interpret. This is because they have fewer variables and relationships to consider. However, simpler models may not be able to capture the complexity of the data as well as more complex models [2]. This can lead to lower accuracy and less reliable predictions.

Balancing these two opposing factors is a significant challenge in the development of XAI systems. Researchers are exploring various approaches to address this challenge. These approaches include developing methods for explaining complex models in simpler terms, using visualization techniques to make models more transparent, and designing models that are inherently interpretable.

By addressing the trade-off between complexity and interpretability, XAI systems can become more powerful and accessible tools for understanding and interacting with artificial intelligence. This will enable us to develop more trustworthy and reliable AI systems that can be used to solve complex problems and make better decisions.

Evaluation of Explanations

Measuring the quality of explanations is a complex and multifaceted task that requires consideration of multiple criteria. Several factors contribute to the effectiveness of an explanation, including completeness, correctness, comprehensibility, engagement, and actionability. Completeness ensures that all relevant aspects and details are covered, while correctness guarantees accuracy and freedom from errors. Comprehensibility makes the explanation easy to understand by presenting it in clear and familiar language. Engagement keeps the reader's attention through storytelling and examples, and actionability provides practical insights that can be applied to real-world situations.

In addition to these criteria, user studies often help evaluate the effectiveness of explanations. These studies measure factors

like understanding, satisfaction, and the likelihood of using the information provided. Developing standardized metrics and evaluation frameworks is crucial for advancing the field and ensuring that explanations are of high quality and meet users' needs [15].

Domain-Specific Requirements

The field of explainable artificial intelligence (XAI) aims to make complex machine learning models more transparent and interpretable to humans. This is particularly important in domains where high-stakes decisions are made based on model predictions, such as healthcare, finance, and criminal justice. However, the requirements for explainability can vary significantly across different domains and user groups.

In the medical domain, for example, healthcare professionals often need detailed causal explanations of why a particular diagnosis or treatment decision was made. This is because medical decisions are often complex and involve a multitude of factors, and doctors need to be able to understand the rationale behind a model's predictions in order to make informed decisions. In contrast, end-users of a recommendation system might only need simple justifications for why a particular item or service was recommended to them. This is because they are typically not interested in the underlying details of the model, but rather just want to know whether the recommendation is trustworthy and relevant to their needs.

Human Factors

Understanding how humans perceive and interact with artificial intelligence (AI) explanations is crucial for designing effective AI systems. Cognitive biases, user expertise, and contextual factors all influence how explanations are received and understood. Research in human-computer interaction (HCI) and cognitive psychology is essential to designing explanations that are truly helpful to users. By studying how people process and evaluate explanations, researchers can create explanations that are clear, concise, and tailored to the user's needs. This can lead to improved trust, understanding, and decision-making when interacting with AI systems.

Future Directions

Integration of Automated Reasoning

To achieve truly explainable AI, integrating automated reasoning with AI models is essential. This involves developing systems that can generate logical explanations based on the input data and the model's decision-making process. Such systems would go beyond merely interpreting or visualizing the model's output, providing coherent, human-understandable explanations.

Integrating causal reasoning into AI models can help provide explanations that not only describe what the model did but also why it made a particular decision. This approach can enhance the transparency and trustworthiness of AI systems by providing deeper insights into their behavior [16].

Human-AI Collaboration

Future XAI systems should focus on enhancing human-AI collaboration by enabling users to interact with and query AI models. This interaction can lead to better understanding and trust, as users can explore the AI's reasoning process and validate its decisions in real-time.

Developing interactive tools that allow users to ask questions about the AI's decisions and receive explanations can improve

transparency and trust. For example, a doctor could query an AI system about why it recommended a particular treatment, receiving detailed explanations that help verify the recommendation.

Ethical and Fair AI

Ensuring that AI systems are ethical and fair is a critical aspect of XAI. Future research should focus on developing techniques to detect and mitigate biases in AI models and providing explanations that highlight potential ethical concerns.

AI systems can be designed to detect and mitigate biases in their decision-making process. Providing explanations that highlight potential biases and show how they were addressed can help ensure that the AI system operates fairly and ethically.

Conclusion

Explainable AI (XAI) aims to make AI systems more transparent and understandable to humans. This is crucial for building trust in AI systems and ensuring that they are used responsibly. There are several approaches to XAI, each with its own strengths and weaknesses. One common approach is to use post-hoc explanation methods, which generate explanations for AI decisions after they have been made. Another approach is to use interpretable models, which are designed to be inherently understandable by humans.

A key concept in XAI is the idea of fidelity. Fidelity refers to the degree to which an explanation reflects the actual behavior of the AI system. It is important for explanations to be both accurate and complete, in order to avoid misleading users. Another important concept in XAI is the idea of counterfactuals. Counterfactuals are hypothetical situations that are similar to the actual situation, but with one or more changes. By examining counterfactuals, we can gain insights into how the AI system would behave in different situations.

Despite the progress that has been made in XAI, there are still a number of challenges that need to be addressed. One challenge is the trade-off between fidelity and simplicity. Explanations that are highly accurate may be difficult to understand, while explanations that are simple may not be very accurate. Another challenge is the need to develop XAI methods that can be applied to a wide range of AI systems. Currently, most XAI methods are only applicable to specific types of AI systems.

Future research should focus on emerging challenges such as model drift, federated learning, and edge computing. Model drift occurs when a model's predictions become less accurate over time due to changes in the underlying data. Federated learning allows collaborative training of models across decentralized devices without sharing sensitive data, providing privacy-preserving ML opportunities. Edge computing brings computation closer to the data source, enabling real-time applications and resource-constrained environments.

References

1. Holzinger Andreas (2018) From Machine Learning to Explainable AI. 2018 World Symposium on Digital Intelligence for Systems and Machines (DISA) 55-66.
2. Rai, Arun (2019) Explainable AI: from black box to glass box. *Journal of the Academy of Marketing Science* 48: 137-141.
3. Hoffman Robert R, Gary Klein, Shane T Mueller (2018) Explaining Explanation For "Explainable AI". *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 62: 197-201.

4. Doran Derek, Sarah Schulz, Tarek R Besold (2017) What Does Explainable AI Really Mean? A New Conceptualization of Perspectives. *ArXiv abs/1710.00794* (2017).
5. Xu Feiyu, Hans Uszkoreit, Yangzhou Du, Wei Fan, Dongyan Zhao, et al. (2019) Explainable AI: A Brief Survey on History, Research Areas, Approaches and Challenges. *Natural Language Processing and Chinese Computing*.
6. Gunning David, Mark Stefik, Jaesik Choi, Timothy Miller, Simone Stumpf and Guang-Zhong Yang. XAI—Explainable artificial intelligence. *Science Robotics* 4.
7. Holzinger, Andreas, Anna Saranti, Christoph Molnar, P. Biecek and Wojciech Samek. "Explainable AI Methods - A Brief Overview." *xxAI@ICML* (2020).
8. Arrieta Alejandro Barredo, Natalia Díaz Rodríguez, Javier Del Ser, Adrien Bennetot, Siham, et al. (2019) Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI. *Inf Fusion* 58: 82-115.
9. Goebel Randy, Ajay Chander, Katharina Holzinger, Freddy Lécué, Zeynep Akata, et al. (2018) Explainable AI: The New 42?. *International Cross-Domain Conference on Machine Learning and Knowledge Extraction* 295–303.
10. Lipton, Zachary Chase (2016) The mythos of model interpretability. *Communications of the ACM* 61: 36-43.
11. Rosenfeld A, Ariella Richardson (2019) Explainability in human-agent systems. *Autonomous Agents and Multi-Agent Systems* 33: 673-705.
12. Lundberg S (2024) An introduction to explainable AI with Shapley values," SHAP. [Online]. Available: https://shap.readthedocs.io/en/latest/example_notebooks/overviews/An%20introduction%20to%20explainable%20AI%20with%20Shapley%20values.html.
13. Hoffman Robert R, Shane T Mueller, Gary Klein, Jordan Litman (2018) Metrics for Explainable AI: Challenges and Prospects. *ArXiv abs/1812.04608* (2018).
14. Fraunhofer Heinrich Hertz Institute (2024) Layer-wise Relevance Propagation," Fraunhofer. [Online]. Available: <https://www.hhi.fraunhofer.de/en/departments/ai/technologies-and-solutions/layer-wise-relevance-propagation.html>.
15. Shin Don Donghee (2021) The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable AI. *Int J Hum Comput Stud* 146: 102551.
16. Gohel Prashant, Priyanka Singh, Manoranjan Mohanty (2021) Explainable AI: current status and future directions. *ArXiv abs/2107.07045*.

Copyright: ©2022 Pushkar Mehendale This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.