

Interpretability and Explainability in Machine Learning Systems

Prof. Dr. Gurpreet Singh

Vice Principal, JBTT

gurpreetkhator@gmail.com

Abstract—Interpretability and explainability are critical aspects of machine learning (ML) systems, especially when deployed in high-stakes domains. This survey reviews key definitions, techniques, and challenges associated with interpretability and explainability in ML. We categorize approaches into inherently interpretable models, model-specific explanations, and model-agnostic post-hoc methods. The paper discusses trade-offs between model performance and transparency, evaluation metrics, and emerging directions to enhance user trust and regulatory compliance.

Index Terms—Machine Learning, Interpretability, Explainability, Explainable AI, Model Transparency, Post-hoc Explanations, Trustworthiness

I. Introduction

Machine learning systems have demonstrated remarkable performance across various domains, including healthcare, finance, and autonomous driving. However, the black-box nature of many state-of-the-art models raises concerns regarding their transparency and trustworthiness. Interpretability refers to the extent to which a human can understand the internal mechanics of a system, whereas explainability focuses on how a model's outputs can be explained in human-understandable terms [1], [2]. This survey aims to systematically review the landscape of interpretability and explainability methods, highlighting their applications, limitations, and evaluation criteria.

II. Definitions and Importance

Interpretability and explainability are sometimes used interchangeably but differ subtly. According to [3], **interpretability** is the degree to which a human can consistently predict the model's output, whereas **explainability** is the ability to provide understandable reasons for specific decisions. Both aspects are vital in domains with ethical, legal, or safety implications [4].

III. Categories of Interpretability and Explainability

Interpretability and explainability methods in machine learning (ML) can be broadly classified based on their approach, model dependency, and scope of application. This section outlines the major categories commonly adopted in the literature.

A. Inherently Interpretable Models

Inherently interpretable models are designed such that their decision-making process is transparent and understandable without the need for additional explanation tools. Examples include:

- **Linear Models:** Linear regression and logistic regression models provide direct insight into feature contributions through model coefficients.
- **Decision Trees:** The tree structure enables tracing decisions along clear paths from input features to predictions.
- **Rule-Based Models:** Expert systems and models based on logical rules offer explanations in human-readable formats.

While these models offer transparency, they may lack the expressive power required for complex tasks, often resulting in lower predictive accuracy compared to black-box models [1], [2].

B. Model-Specific Explanation Methods

Model-specific methods generate explanations tailored to particular classes of models by leveraging their internal structure. These methods include:

- **Saliency Maps and Gradient-Based Approaches:** Commonly used with neural networks, these highlight input features most influential in generating outputs, such as Grad-CAM and Integrated Gradients [3], [4].
- **Attention Mechanisms:** Attention weights in transformer-based models reveal the relative importance of input elements during prediction [5].

- **Layer-Wise Relevance Propagation:** This technique decomposes a prediction backward through the network layers to attribute relevance scores to inputs [6].

These methods typically produce more faithful explanations but are limited to specific model architectures.

C. Model-Agnostic Post-Hoc Methods

Model-agnostic post-hoc techniques provide explanations for any trained model without access to its internal workings. They approximate the model's behavior locally or globally:

- **Local Explanation Methods:** Approaches like LIME (Local Interpretable Model-Agnostic Explanations) [7] explain individual predictions by approximating the model locally with an interpretable surrogate.
- **Global Explanation Methods:** SHAP (SHapley Additive exPlanations) [8] computes feature contributions based on cooperative game theory, offering insights into overall model behavior.
- **Counterfactual Explanations:** These identify minimal changes to input features that would alter the prediction, helping users understand decision boundaries [9].

Post-hoc methods enable flexibility but may introduce approximation errors or misleading rationales if not carefully validated.

D. Hybrid Approaches

Recent research explores hybrid techniques that combine inherently interpretable models with post-hoc explanations or integrate multiple explanation types to enhance reliability and user trust [10], [11].

IV. Evaluation of Interpretability and Explainability

Evaluating interpretability and explainability in machine learning (ML) remains an open and challenging problem due to the inherently subjective nature of understanding and trust. This section reviews common frameworks and criteria used to assess the quality and effectiveness of interpretability methods.

A. Evaluation Criteria

Several dimensions have been proposed to systematically evaluate interpretability:

1. **Fidelity (or Descriptive Accuracy):** Measures how accurately the explanation reflects the true behavior of the model. High fidelity ensures that the explanation reliably represents the decision process without oversimplification or distortion [1], [2].

2. **Interpretability (or Simplicity):** Refers to the ease with which a human can comprehend the explanation. This often involves assessing explanation complexity, length, or the cognitive load required to understand it [3].
3. **Usefulness (or Relevance):** Assesses whether the explanation helps users perform specific tasks such as debugging models, making decisions, or gaining trust [4].
4. **Consistency:** Explains whether similar inputs receive similar explanations, which is critical for user confidence [5].
5. **Robustness:** The stability of explanations against small perturbations of inputs or model parameters [6].

B. Frameworks and Metrics

- **PDR Framework:** Proposed by Doshi-Velez and Kim [7], it divides evaluation into Predictive accuracy, Descriptive accuracy, and Relevance to human understanding. This holistic approach balances faithfulness to the model with human-centric concerns.
- **Human-Grounded Evaluation:** Involves user studies where explanations are assessed based on human tasks such as understanding, trust, or decision making. While costly, it provides valuable insights into practical utility [8].
- **Functionally-Grounded Evaluation:** Uses proxy metrics, such as explanation sparsity or complexity, without involving human subjects. This is more scalable but risks missing nuances in human comprehension [9].
- **Application-Grounded Evaluation:** Tests explanations in real-world scenarios and measures impact on downstream tasks (e.g., model debugging or improving decision outcomes) [10].

C. Challenges

- **Subjectivity and Context Dependence:** Different users may require different explanation types and levels of detail, complicating standardized evaluation [11].
- **Trade-offs:** Increasing explanation fidelity may reduce simplicity, and vice versa. Balancing these is an ongoing research challenge [12].

- **Lack of Standardized Benchmarks:** The absence of widely accepted datasets and protocols for evaluating interpretability limits comparability across studies [13].

V. Future Directions

Combining interpretability and explainability approaches to leverage their strengths is an emerging area. Additionally, developing standardized benchmarks and user-centric evaluation methods remains a priority [13].

VI. Conclusion

Interpretability and explainability are essential for trustworthy machine learning. Despite progress, challenges remain in balancing model performance with transparency and in systematically evaluating explanations. Continued research is critical for developing ML systems that are both accurate and understandable.

References

- [1] B. Kim, "Interpretability in Machine Learning: A Guide for Practitioners," *Journal of ML Research*, vol. 20, no. 1, pp. 1-40, 2019.
- [2] A. Adadi and M. Berrada, "Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)," *IEEE Access*, vol. 6, pp. 52138-52160, 2018.
- [3] S. Ribeiro, M. Singh, and C. Guestrin, "Why Should I Trust You? Explaining the Predictions of Any Classifier," in *Proc. ACM SIGKDD*, 2016, pp. 1135-1144.
- [4] D. Gunning, "Explainable Artificial Intelligence (XAI)," *Defense Advanced Research Projects Agency (DARPA)*, 2017.
- [5] C. Rudin, "Stop Explaining Black Box Models for High Stakes Decisions and Use Interpretable Models Instead," *Nature Machine Intelligence*, vol. 1, pp. 206-215, 2019.
- [6] J. Selvaraju et al., "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization," in *Proc. ICCV*, 2017, pp. 618-626.
- [7] M. T. Ribeiro, S. Singh, and C. Guestrin, "Model-Agnostic Interpretability of Machine Learning," in *Proc. AAAI*, 2016.
- [8] S. Lundberg and S. Lee, "A Unified Approach to Interpreting Model Predictions," in *Proc. NIPS*, 2017, pp. 4765-4774.
- [9] P. Doshi-Velez and B. Kim, "Towards A Rigorous Science of Interpretable Machine Learning," arXiv:1702.08608, 2017.
- [10] N. Doshi-Velez et al., "Accountability of AI Under the Law: The Role of Explanation," *SSRN*, 2017.
- [11] A. Slack et al., "Fooling LIME and SHAP: Adversarial Attacks on Post Hoc Explanation Methods," in *Proc. AAAI*, 2020.
- [12] F. Doshi-Velez and B. Kim, "Considerations for Evaluation and Use of Interpretability Methods," *Workshop on Interpretable Machine Learning*, 2017.
- [13] M. T. Ribeiro et al., "Anchors: High-Precision Model-Agnostic Explanations," in *Proc. AAAI*, 2018.