

UNLOCKING THE PATH FORWARD: NAVIGATING CHALLENGES & EMBRACING OPPORTUNITIES IN EXPLAINABLE ARTIFICIAL INTELLIGENCE (XAI)

Abstract

With the continuous digitization and proliferation of digital technologies and platforms, we are embarking on a new age of Artificial Intelligence applications considering machine learning and deep learning as the core technologies. AI systems leveraging machine learning models, have a profound impact on how Business and IT leaders make real-time critical decisions without enough context on how these complex systems arrive to a particular outcome. Essentially, they simply rely on black-box technology which erode their trust on AI systems due to lack of transparency, fairness and the ability to explain the outcome of machine learning models. Our study focuses on the importance of ML Models and underlying algorithms assessment in terms of the explanations they produce, their constraints, and real-world uses. Through a meticulous analysis of their advantages and limitations, the objective of this paper is to highlight how XAI techniques can potentially reconcile the divergence between the accuracy of machine learning models and the dire need for comprehensible interpretations. This paper specifically talks about the current XAI challenges and future opportunity to design AI Systems which are ethical in the purpose emphasizing the desire for transparency, fairness and interpretability.

Keywords: trustworthy AI; Transparent AI, Explainability, Explainable AI, XAI etc.

Authors

Mr. Jitendra Maan

Head, AI & Cognitive Experience

Software & Services Unit

Tata Consultancy Services

Gurgaon, Haryana, India

Jitendra.maan@tcs.com

I. INTRODUCTION

Explainable AI as a field is constantly evolving as new challenges and opportunities arise due to significant advancement in AI and machine learning technologies. Apparently, Explainable AI becomes more critical given the fact that most organizations are exploring how they can leverage transformer architecture based Large Language Models which are incredibly powerful with remarkable performance across various natural language tasks. Henceforth, Trustworthy Explainable AI plays a vital role in addressing such global problems which is not as easy as it sounds.

The current generation of AI systems offer tremendous benefits, but their effectiveness will be limited by the machine's inability to explain its decisions and actions to users. Explainable AI will help to understand the rationale for the system's decisions to understand, appropriately trust, and effectively manage the AI system. Any XAI system design will consider the following factors –

1. Should provide an explanation of individual decisions made by the ML models.
2. Understanding of the strengths & weaknesses
3. Provide the context on how the system will behave in the future.
4. Strategy to correct the system's mistakes/errors.

The most innovative approach is to create a suite of machine learning techniques that produce more explainable models, while maintaining high accuracy and performance.

II. LITERATURE SURVEY

Several XAI techniques have been proposed by different researchers from time to time, certainly, SHAP (SHapley Additive exPlanations) is indeed a significant method in the field of Explainable AI (XAI) wherenoteworthy prior works and research is carried out with the aim to provide a unified framework for interpreting machine learning models.

Topol's work [1] for the convergence of human and artificial intelligence explores the intersection of human and AI capabilities in healthcare sector and author has emphasizes the transformative impact of AI in enhancing medical data analysis, personalized treatments, and patient outcomes.

Montavon, Samek, and Muller [2] delve into techniques for elucidating the inner workings of deep neural networks (DNNs) and their study explores various methods to unravel the decision-making processes of DNNs, aiming to enhance transparency and interpretability. Here, authors presented techniques like visualization, attribution, and sensitivity analysis, which shed light on feature importance and model behavior. Guidotti R, Monreale A and Ruggieri S [3] conducted an extensive exploration of techniques designed to elucidate the functioning of black box models. Their survey provides a comprehensive resource for researchers and practitioners seeking to enhance interpretability in black box models and recommended various strategies like perturbation- based methods, rule extraction, and gradient-based approach.

European Union's groundbreaking General Data Protection Regulation (2018) [5] on the everyday deployment of machine learning algorithms made a significant impact by imposing constraints on the utilization of automated individual decision-making, pertinently encompassing algorithms that solely take critical decisions biased to developer understanding and interpretations. Furthermore, it will effectively instate a "right to elucidation," empowering users to solicit an explication regarding algorithmic decisions that have substantially impacted them. Adadi and Berrada [6] undertake a comprehensive investigation into the realm of Explainable AI (XAI). The study, published in IEEE Access, systematically explores a range of XAI techniques designed to enhance the interpretability of black-box AI models. The authors provide a structured overview of various methods, including rule-based approaches, post-hoc explanations, and model-specific techniques.

Lipton Z.C. [7] in his paper, challenged the prevailing assumptions about the true nature and achievability of interpretability in complex machine learning models. The author highlighted that the pursuit of perfect transparency might not always align with the complexity of modern models, and interpretable features can often be subjective. His paper stimulates a thought-provoking discourse on the limitations and practicalities of model interpretability, urging a nuanced and context-aware approach to the subject. On the other end, Angwin, J., Larson, J., Mattu, S. and Kirchner, L [8] explored the case of a risk assessment tool employed to predict future criminal behavior and assess its fairness across different racial groups. They reveal that the algorithm disproportionately labeled black defendants as high risk even when compared to white defendants with similar characteristics. Their research paper shed light on the unintended biases inherent in some AI systems and underscores the importance of addressing algorithmic fairness and accountability in critical domains like criminal justice.

III. KEY CHALLENGES IN EXPLAINABLE AI

It's important to understand that the trade-off between model performance and explainability is not always absolute. Ongoing research in the field of XAI aims to mitigate this trade-off by developing hybrid approaches that attempt to retain both high performance and meaningful explanations. The challenge lies in striking the right balance between these two essential aspects, ensuring that AI models are not only accurate but also transparent and trustworthy. There are many more challenges as depicted below:

- 1. Complexity of AI Models:** Many AI models, such as deep neural networks, Generative Adversarial Networks (GANs) are highly complex and opaque, kind of black-box in nature. Understanding their decision-making processes and how they arrive at specific predictions is quite challenging. Explainable AI methods often involve simplifying the model's internal workings to generate human-readable explanations. This simplification can lead to a reduction in the model's complexity, potentially sacrificing the ability to capture intricate patterns and nuances in the data.

When designing AI systems, developers often need to choose between highly accurate yet complex models and simpler models that are easier to explain. This trade-off involves determining the right balance between accuracy and interpretability based on the application's requirements.

2. **Trade-Off between Performance and Explainability:** Explainable AI (XAI) techniques often create a trade-off between performance and explainability due to the inherent complexity of many advanced machine learning models. These techniques aim to make AI models more transparent and interpretable, but achieving this transparency can come at the cost of model performance, whereas simpler models tend to be more interpretable, but they may sacrifice predictive power compared to more complex models.
3. **Model Distillation:** There may be quite a few situations where complex ML models are distilled into simpler ones for the sake of explanation, but the distilled model might not fully capture the original model's performance. Distillation inherently involves a compression of information, leading to a performance drop.

There are several XAI methods that rely on simplified feature representations or surrogate models. These simplifications might discard some of the original model's complexity, impacting its performance.

4. **Context Sensitivity:** Explanations can vary depending on the context and the user's background knowledge. Finding a universally understandable way to present explanations is quite difficult. Besides this, in certain AI systems, models are continuously updated and improved based on new data. This can lead to changing explanations over time, making it hard for users to trust the model's stability and consistency.
5. **Local v/s Global Interpretations:** Some XAI methods provide explanations for individual predictions (local interpretability), while others aim to explain the model's behavior across the entire dataset (global interpretability). Achieving both simultaneously can be challenging and might impact overall performance.
6. **Lack of Standards and Metrics:** The absence of standardized evaluation metrics for explainable AI methods makes it difficult to objectively compare different approaches. The ultimate goal of XAI is to enhance user understanding of AI decisions.

However, there's a lack of standardised metrics to gauge how well users comprehend the provided explanations. Suppose two XAI methods are employed to explain loan approval decisions. While one method uses technical jargon, the other presents explanations in plain language. Without standardised metrics, it's hard to measure which method truly enhances user understanding.

7. **Model Agnostic Evaluation Metrics:** XAI techniques often need to work across different types of machine learning models. The absence of standardized model-agnostic metrics makes it challenging to determine if an explanation method performs consistently across various model architectures. Say for example, A model-agnostic explanation method is applied to explain decisions from both neural networks and decision trees. Without standardized metrics, it's challenging to objectively compare the performance and reliability of the explanations on different model types.
8. **Regularity vs. Novelty Detection:** Highly accurate models might excel at detecting regular patterns in data but struggle with identifying novel, out-of-distribution examples. Explanation methods that focus on regular patterns might overlook unusual cases.

IV. DATA PRIVACY AND SECURITY

Some explainable AI methods rely on inspecting the model's internal workings, which could raise concerns about data privacy and potential vulnerabilities to attacks. XAI methods often involve providing explanations by analysing the internal workings of AI models. If not designed carefully, these explanations might inadvertently reveal sensitive information about individuals or entities in the training data. Say for example, lets imagine an XAI technique used to explain a credit scoring model. If the explanation reveals specific personal attributes or financial details of individuals, it could potentially breach their privacy.

Such examples underscore the complex interplay between XAI, data privacy, and security. While XAI techniques are designed to provide transparency and insights into AI model behavior, they must be developed with robust privacy-preserving mechanisms to avoid inadvertently disclosing sensitive information or being susceptible to adversarial attacks. Addressing these challenges is crucial to ensuring that XAI methods contribute positively to both model transparency and user privacy. Molnar, C., Casalicchio, G., Bischl, B.[10] in their paper explored diverse methods for enhancing model transparency and explainability, addressing the growing need for AI systems to be more understandable and accountable. Their piece of work serves as a valuable resource, offering insights into the landscape of IML techniques and the ongoing efforts to make AI more interpretable across various domains.

V. WHY XAI IS IMPORTANT?

- 1. Cultivating Confidence in AI Models with Explainable AI (XAI):** Let me consider a scenario where a medical professional relies on an AI algorithm to detect a critical medical condition. In such life-and-death situations, the precision of the diagnosis is paramount, underscoring the doctor's reliance on the AI model's competence in decision-making. Explainable AI (XAI) emerges as a pivotal mechanism for nurturing such trust on AI Systems. By furnishing elucidations that enable the doctor to comprehend the rationale behind the AI model's diagnosis, XAI engenders a sense of confidence in the AI system's judgments.
- 2. Building Trust and Establish Transparency:** As AI and machine learning algorithms make critical decisions that impact various aspects of our lives, from financial transactions to medical diagnoses, there's a growing need for transparency in these decisions.
- 3. Avoiding AI Models to be Seen as Inscrutable Black Boxes:** Unlike traditional software where the source code is visible, AI models are viewed as "black boxes" since they use proprietary algorithms which are highly complex. With XAI, users can see inside the "black boxes" and see how the models arrived at their decisions.
- 4. Building a Sustainable and Responsible Ecosystem:** XAI promotes transparency, fairness, and interpretability, ensuring that AI models are inclusive, accountable, and transparent. This will build a sustainable ecosystem that instills trust in end-users.

VI. XAI IN THE CONTEXT OF LARGE LANGUAGE MODELS (LLMs)

AI, Natural Language Processing (NLP) and LLM-based model are complex and most often considered as “Black-box” models and at times, it becomes challenging to understand the specific decisions made by these complex LLM models which typically used to do more complex tasks. Having said, XAI is gaining significant business traction in the context of large language models. There are several factors to consider –

- 1. Trust and Interpretability:** Every industry has one or the other use cases of Large Language Models and non-deterministic behavior of Large Language Models make it a daunting task to interpret how these models yield different outputs for identical inputs as they fit for complex reasoning situations. To build trustworthy relationships with customers, partners and vendors, it is essentially important to provide explanation for the decision arrived by machine learning models. Henceforth, the need of the hour is to mitigate challenges in making AI trustworthy, fair, reliable and explainable with human in the loop.
- 2. Ethical and Regulatory Compliance:** Most of the regulatory bodies are in its nascent stage of developing specifications and standards to ensure transparency and accountability in AI decision-making. Explainable AI works as a catalyst in satisfying the regulatory requirements with the significant focus on ethical and responsible AI.
- 3. Biases and Fairness:** It is important to evaluate the fairness of the model against sensitive attributes (Age, Race, Gender etc) and identify the innovative solutions to mitigate the biases leveraging Generative Adversarial Network (GAN) based Modeling approach. There are lot of instances where Large Language Models get learned on biased patterns present in the training dataset. Here, XAI techniques can help in addressing such biases while making model more fair and unbiased.
- 4. Human-Machine Collaboration:** In many business scenarios, LLM based AI Models not only factor into the human decision-making but it can foster a more collaborative approach where man and machine work together and complement each other on specific activities and there is a dire need to establish/enable a more collaborative approach to provide a more conducive environment for both human and machine to work effectively.
- 5. Explainable AI to Aid Debugging and Improvement:** Due to complexity, Large Language Models are quite prone to mistakes which may lead to incorrect predictions which opens up the opportunity for developers to debug and improve the AI system. Explainable AI can provide such insights on why a particular model generates an output, which helps to understand how the model works with the focus on improving and promoting responsible AI usage with the fast-changing context or environments.
- 6. Hallucination:** It is quite a pertinent problem in most of the Large Language Models. One of the most common reason for such problems is due to the model which is overfit on either biased dataset or dataset is too small which generates output which is not a true representation of real world scenarios. Explainable AI can help to place guardrails to avoid unintended consequences.

- 7. Ease of Model Selection:** There are lot of B2C applications where AI system is directly interface with consumers/end-users (say for example Intelligent Virtual Assistants and Chatbots) and transparent explanation of model decisions can help to select the most appropriate model for a specific task in the given context.

VII. EXPLAINABLE AI IMPACT ON SOCIETY

- 1. Transparency and Accountability:** XAI promotes transparency by enabling users and stakeholders to understand the decision-making processes of AI systems. This accountability is crucial in domains like finance, healthcare, and criminal justice, where decisions have significant real-world consequences. XAI helps prevent biased or unjust decisions from being obscured by "black-box" models.
- 2. Ethical AI Deployment:** XAI plays a pivotal role in ensuring the ethical deployment of AI technologies. It allows organizations to identify and rectify biases, discriminatory patterns, and unfair outcomes, thus reducing the risk of perpetuating societal inequalities through AI systems.
- 3. Trust and Adoption:** The interpretability provided by XAI fosters trust between users and AI systems. Users are more likely to embrace AI technologies if they understand how decisions are made, leading to higher adoption rates across various industries and applications.
- 4. User Empowerment:** XAI empowers users to make informed decisions by understanding AI-generated recommendations. This is particularly relevant in healthcare, where patients can better comprehend medical diagnoses and treatments, leading to shared decision-making between doctors and patients.
- 5. Regulatory Compliance:** As regulatory bodies worldwide focus on AI ethics and accountability, XAI helps organizations comply with regulations by providing explanations and justifications for AI-driven decisions. It aids in ensuring that AI applications adhere to legal and ethical standards.
- 6. Education and AI Literacy:** XAI techniques are valuable educational tools, enhancing public understanding of AI principles. Society becomes more AI-literate, enabling individuals to critically assess and navigate AI-generated information and outcomes.
- 7. Fairness and Bias Mitigation:** XAI helps identify biases and discrimination in AI models, enabling mitigation strategies. This leads to the development of fairer and more equitable AI systems that do not discriminate against individuals based on their characteristics.
- 8. Human-AI Collaboration:** In collaborative settings, such as medical diagnosis or financial advising, XAI enables effective collaboration between humans and AI. Users can trust AI's suggestions and provide meaningful feedback, creating a harmonious human-AI partnership.

- 9. Crisis Management and Decision Support:** During crises like disease outbreaks or natural disasters, XAI can provide insights into the rationale behind AI-driven decisions. This assists policymakers in making informed choices to mitigate the impact of such events.
- 10. Policy and Regulation Development:** Policymakers and regulators benefit from XAI insights to create policies and regulations that govern AI systems effectively. XAI research helps inform the development of standards that balance innovation and societal well-being.

In essence, the impact of XAI on society is multidimensional, encompassing improved transparency, ethical AI deployment, enhanced trust, and more informed decision-making. It contributes to a responsible AI ecosystem that aligns AI advancements with societal values and concerns.

VIII. EXPLAINABLE AI: FUTURE OPPORTUNITIES

There are significant areas where Explainable AI can address gaps and provide future opportunities.

- 1. Legal and Ethical Considerations:** Ethical and legal concerns are rapidly growing into AI Systems and machine learning models due to biased datasets with very little visibility on explanation of how machine learning model predicted the outcome. XAI system should be designed considering legal and ethical concerns to address potential biases and discrimination.
- 2. Explainable Reinforcement Learning (XRP):** XRP opens up new door for AI Systems to explain sequential decision making. Creating a self-driving car AI that not only navigates traffic but also offers clear explanations for its decisions, enhancing passenger trust and safety.
- 3. Guidelines and Standardization:** XAI solutions and innovative techniques will help research community and industry practitioners in selecting the most viable methods based on standardized metrics/KPIs, performance checklist and guidelines.
- 4. Collaborative Human-AI Decision Making:** Here the opportunity is to establish the frameworks where humans and AI systems collaboratively make decisions, facilitated by clear and meaningful AI explanations. Say for example, creating a financial investment platform that provides real-time explanations for AI-generated investment suggestions, allowing users to work in tandem with the AI for more informed decisions.
- 5. Interpretable AI Model Architectures:** There is always a trade-off between model performance and explainability. Future solutions will strike a balance between them by deploying interpretable AI models contextualized to Industry domains.
- 6. Explainable Deep Learning:** Key opportunity is to develop more effective methods for explaining the decisions of deep learning models, which are known for their complexity and black-box nature. Say for example, enhancing an image recognition AI's ability to

explain why it classified a certain image as a specific object, revealing the specific features that influenced it

IX. CONCLUSIONS

Despite these challenges, the future of XAI holds a promise of transformative opportunities. As algorithms mature, the field will likely witness the emergence of refined methodologies that harmonize the intricacies of complex models with user-friendly explanations. Enhanced transparency will cultivate trust, encouraging broader acceptance of AI's role in shaping various domains. The growth of standardized evaluation metrics will offer a solid foundation for benchmarking and assessing different XAI techniques, advancing the field in a systematic manner. Moreover, the synergy between XAI and fairness will create AI systems that are not only accurate but also free from biases, safeguarding against discrimination and ensuring equitable outcomes. Collaborative human- AI decision-making, fostered by comprehensible explanations, will drive innovative solutions across sectors, empowering users to make informed choices in partnership with intelligent machines.

In this ongoing journey, XAI will continue to empower individuals with AI literacy, propelling society towards an era where AI is embraced not as an enigma, but as a transparent and accountable tool. As researchers, practitioners, and policymakers engage in the pursuit of XAI's potential, they embark on a transformative expedition, one that empowers AI systems to not just predict, but to explain, inspire, and positively impact the world in a manner that resonates with the principles of transparency, fairness, and collaboration.

REFERENCES

- [1] Topol, E.J. High-performance medicine: The convergence of human and artificial intelligence. *Nat. Med.* 2019, 25, 44–56.
- [2] Montavon G, Samek W, Iler KRM. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*. 2018;73:1–15.
- [3] Guidotti R, Monreale A, Ruggieri S, et al. A Survey of Methods for Explaining Black Box Models. *ACM Comput Surv.* 2018 August;51(5).
- [4] Adadi A, Berrada M. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access.* 2018;6:52138– 52160.
- [5] Goodman, B.; Flaxman, S. European Union regulations on algorithmic decision-making and a “right to explanation”. *AI Mag.* 2017, 38, 50–57
- [6] Adadi, A.; Berrada, M. Peeking inside the black-box: A survey on Explainable Artificial Intelligence (XAI). *IEEE Access* 2018, 6, 52138–52160.
- [7] Lipton, Z.C. The mythos of model interpretability. *Queue* 2018, 16, 31–57.
- [8] Angwin, J.; Larson, J.; Mattu, S.; Kirchner, L. Machine bias. *ProPublica* May 2016, 23, 139–15
- [9] Imana, B.; Korolova, A.; Heidemann, J. Auditing for Discrimination in Algorithms Delivering Job Ads. *arXiv* 2021, arXiv:2104.04502
- [10] Molnar, C.; Casalicchio, G.; Bischl, B. Interpretable Machine Learning—A Brief History, State-of-the-Art and Challenges. *arXiv* 2020, arXiv:2010.09337
- [11] Birhane, A. Algorithmic injustice: A relational ethics approach. *Patterns* 2021, 2, 100205