

ILLUMINATING THE BLACK BOX: A DEEP DIVE INTO SHAP FOR EXPLAINABLE AI

Piyush Zaverbhai Patel¹, Dr. Vipul Vekariya², Dr. Kruti Sutaria³

¹M.tech Student, Parul Institute of Engineering and Technology, Vadodara

²Dean, Parul Institute of Engineering and Technology, Vadodara

³Professor, Parul Institute of Engineering and Technology, College University Name: parul University, Vadodara

piyush.pp1096@gmail.com¹
enggdean@paruluniversity.ac.in²
kruti.sutaria25509@paruluniversity.ac.in³

Abstract

The fast development of artificial intelligence (AI) has produced complex models that are commonly viewed as black boxes, which restrict transparency and interpretability. Based on cooperative game theory, SHAP (SHapley Additive exPlanations) has proven to be an effective solution to this problem, contributing to feature values consistently and in an interpretable way. This review elaborates on the theoretical basis, uses, and difficulties of SHAP in explainable artificial intelligence (XAI). Our fields of analysis include literature in healthcare, natural language processing, finance and customer behavior analytics. The review mentions strengths of SHAP in generating locally and globally consistent explanations, as well as dealing with computational complexity and high-dimensional data. Scalability, domain specific knowledge integration and user-centered evaluation metrics are identified gaps in research. This review will present a summary of the role of SHAP in promoting trustful AI by synthesizing recent research.

Keywords: Explainable AI, SHAP, interpretability, model transparency, black box models, feature attribution

Introduction

Background and Context

Artificial intelligence has been a success in areas including healthcare diagnostics, finance, marketing, and natural language processing. There are however numerous complex models, especially deep learning models, that are opaque and do not generally give an understanding of how the predictions are conducted (Oxley, 2020). Such inability to interpret has created issues of trust, fairness, responsibility, and compliance with regulations. Explainable AI (XAI) has emerged as an important area of study to solve these problems (Tandon et al., 2024).

SHAP, the model created by Lundberg and Lee (2017) is one of the most theoretically based interpretability tools. SHAP is a single, additive example of feature attribution based on Shapley values of cooperative game theory, with both local and global descriptions of the model.

Importance and Significance of the Review

SHAP is important because it has the potential to offer robust explanations improving transparency in a model, making stakeholders more trustful. It is used in a wide variety of applications: to enhance clinical decision-making in healthcare (Muhammad et al., 2024) or predicting consumer behavior (IJSRA, 2025). As the importance of ethical and trustworthy AI continues to grow, an overview of the strengths, limitations, and future projections of SHAP is critical to researchers and practitioners alike.

Objectives and Scope

This review aims to:

1. Explore SHAP's theoretical underpinnings and methodological approaches.
2. Examine its applications across multiple domains.

3. Analyze comparative advantages and challenges against other explainability techniques such as LIME.
4. Identify research gaps and propose future research directions.

The scope of this review focuses on SHAP and its variants, excluding in-depth technical details of unrelated XAI methods unless used for comparison.

Research Gaps Identified

Despite its growing use, SHAP faces challenges in computational scalability, usability for non-expert audiences, and adaptability to highly dynamic data. Few studies have systematically assessed SHAP's effectiveness in real-world deployments, leaving a gap between theory and practice (PMCID: PMC10997568, 2024). Additionally, user-centric evaluation frameworks for explanation quality remain underdeveloped.

Methodology

This review adopts a **systematic literature review approach** designed to capture, analyze, and synthesize the most relevant studies concerning SHAP (SHapley Additive exPlanations) within the broader field of explainable AI (XAI). The methodology was carefully structured to ensure transparency, replicability, and relevance.

First, multiple academic databases were consulted, including **Scopus, Web of Science, IEEE Xplore, and Google Scholar**, as these are among the most comprehensive sources for peer-reviewed AI and computer science literature. To supplement academic articles, reputable industry blogs and technical reports were also considered, given the practical nature and rapid adoption of SHAP in real-world applications.

A combination of keywords was employed to identify the literature: "*SHAP*," "*explainable AI*," "*SHapley additive explanations*," "*model interpretability*," "*black box AI*," and "*XAI applications*." Boolean operators were used where possible to expand or narrow search results. The search was restricted to studies published between **2017 and 2025**, corresponding to the period following the introduction of SHAP by Lundberg and Lee (2017).

There were definite inclusion and exclusion criteria. Inclusion criteria were; (1) peer-reviewed journal articles or conference proceedings, (2) systematic reviews or narrative reviews on SHAP, and (3) reputable industry or technical reports that contained empirical data regarding the implementation of SHAP. Filtering criteria eliminated: (1) articles not clearly methodologically rigorous, (2) articles only comparing irrelevant interpretability frameworks to SHAP, and (3) not written in English.

The first search resulted in 43 articles that were relevant. They were filtered in two phases: Title/abstract screening and a full-text review. Through this process, 15 core sources have been identified to be analyzed in depth according to their direct relevance, methodological soundness and contribution to the understanding of the theoretical foundations, applications and limitations of SHAP.

This methodology approach made sure that not only did the review include foundational theoretical works but also included practical domain specific applications of SHAP to explainable AI, which is a balanced viewpoint of SHAP in explainable AI.

Literature Review

1. Theoretical Foundations of SHAP

Lundberg and Lee (2017) put the background of SHAP in perspective by generalizing the concept of cooperative game theory, namely Shapley values, to machine learning interpretability. Values of the game that were originally designed to distribute the payoffs of a cooperative game in a fair

way were redefined to distribute the contribution of model prediction among the features of input. The adaptation results in SHAP, a general additive model which provides uniform and locally correct explanations. The authors demonstrated that the SHAP explanations satisfy three desirable properties, i.e., local accuracy, missingness, and consistency. Local accuracy: Values of the explanation sum to the model output, missingness: Relevant features are not assigned any value, and consistency: SHAP values of a feature will not decrease when the model is perturbed to make the feature have a bigger contribution to the model output. These conceptual strengths give SHAP an edge over earlier interpretability methods, which were heuristic.

Lundberg, Fu, and Lee (2020) also explained that SHAP is theoretically sound compared to alternative methods, including LIME. They also pointed out the additive properties of SHAP that enable the explanations to be mathematically consistent and to be interpreted using different model types. In contrast to feature importance scores, which do not always consider the local context, SHAP is able to directly attribute predictions to each prediction, allowing its explanation to be globally and locally explainable. This has made SHAP an important tool in areas where accuracy is not a crucial factor, but accountability and legibility is crucial as in medicine and finance.

According to Oxley (2020), the conceptual beauty of SHAP is that it addresses the so-called black box problem of machine learning. Not only does SHAP decompose predictions into feature contributions, but it also provides information on the relative significance of variables within a set of predictions, as per Oxley. The linkage of interpretability at an instance level to the general behavior of a model has been a defining feature of this property, and contributed to the encouragement of its usage by researchers and practitioners.

2. SHAP vs. Other Interpretability Methods (LIME and Beyond)

Lundberg, Fu, and Lee (2020) has turned out to be one of the most impactful comparative studies. They claimed that; unlike the two approaches that both aim at estimating the relative importance of features in a prediction, LIME uses local surrogate models to produce an explanation. This surrogate-based defence is also volatile and inconsistent, that is, a slightly different perturbation will still provide a different explanation of the same prediction. By comparison, the uniform values of Shapley are the explanations in SHAP that are based on a theory of explanation. This renders SHAP more costly to compute and, most notably, much more reliable, at least in high stakes situations, such as in healthcare or criminal justice.

Arize AI (2024) identified this feature and placed SHAP alongside the best tools in the interpretability space. The report has noted that LIME is still a viable alternative with respect to computing efficiency and convenience, and SHAP is preferable to LIME with respect to accuracy, reproducibility, and reliability. The authors further added that SHAP descriptions are additive such that they may be added together to form global descriptions of model behavior which cannot be similarly done in a LIME.

Contreras et al. (2024) presented novel ideas to improve the interpretability capabilities of SHAP and LIME and suggested spectral zone-based methods. Their paper showed that SHAP and LIME explanations can be easier to understand by grouping perturbations, particularly in the context of analytical chemistry. Although the study did not reject the utility of LIME, it pointed out SHAP being more flexible in perturbation of groups due to its additive and consistent base. This shows that although LIME is very popular as a benchmark, SHAP has become the standard metric against which new interpretability methods are evaluated.

Opensense Labs (2025) showing why organizations choose SHAP to complete trustful AI systems instead of LIME. They said that the stakeholders usually need explanations which are not only

mathematically hearing but also practical and simple to communicate. Whereas insights in seconds might be enough in prototyping, enterprises favor SHAP in production-level system where interpretability needs to be audited and justified to regulators. The superiority of SHAP over LIME is therefore not only hypothetical but practical in both regulatory and ethical terms.

3. SHAP in Healthcare and Medicine

Muhammad et al. (2024) highlighted the revolutionary nature of SHAP in the analysis of medical imagery. They emphasized in their systematic review that SHAP has the potential to demystify the use of convolutional neural networks to detect tumors, allowing clinicians to understand what pixels or regions contributed to their predictions. SHAP enhances transparency beyond what feature attributions can provide and assists in clinical decision-making by confirming that models are predicting, not spurious relationships, based on medically meaningful signals. The authors have concluded that the implementation of SHAP in healthcare can greatly enhance trust and enable the collaboration between AI systems and human specialists.

The systematic review by **PMC10997568 (2024)** on SHAP-based explainability in healthcare revealed several main benefits, namely, discovery of aspects of model bias and the possibility to provide patient-level information. As mentioned in the review, SHAP explanations can be used to determine when a diagnostic model is unreasonably biased or unreasonably penal against a particular group of patients, thus overcoming the issue of unfairness in clinical practice. Nevertheless, the authors also admitted that computational difficulties were present, particularly when working with large-dimensional genomic and imaging data.

Tandon et al. (2024) talked about the more general application of SHAP in explainable AI, and healthcare is one of the most necessary fields. They described how clinicians tend to be suspicious of black-box AI models, especially in situations where the safety of patients is involved. SHAP fills this gap by providing clear explanations which not only explain the contribution of features but also follow domain knowledge. To illustrate this point, in cardiovascular risk prediction, SHAP explanations typically give prominence to medically verified risk factors like age, cholesterol levels, and smoking history, providing credibility to the outputs of AI.

Muhammad et al. (2024, Journal of Healthcare Informatics Research) extended the scope of discussion to other healthcare applications, including treatment outcome predictions or drug response modeling. They observed that SHAP explanations are also being applied more in analyzing electronic health records (EHR) to determine risk factors when a patient is readmission. SHAP is applied in this type of application to increase the comprehension of the complex machine learning models presently used by clinicians in ranking the relevant features and giving each patient a personalized explanation.

The healthcare utility of SHAP was expanded to predicting behavior in the field of public health by **International Journal of Science and Research Archive (2025)**. They demonstrated that SHAP-based models could be used to make interpretable predictions of patient adherence to treatment regimens to inform policymakers in designing interventions specific to specific populations. This focus on interpretability means that AI-driven insights will allow trusting the insights without depending on opaque decision-making processes by the officials in the sphere of public health.

4. SHAP in Finance and Risk Management

Praxis AI (2025) highlighted, financial institutions are progressively relying on SHAP to comply with the regulatory requirements of transparency in algorithmic decision-making. Complex models are frequently employed to credit score, detect fraud, and assess the risk of financial

systems, with explainability being a regulatory requirement (as under the GDPR and Fair Credit Reporting Act). SHAP also enables banks and insurers to explain in feature-by-feature terms why they approved or rejected a loan application, showing a customer and an auditor why a loan was approved or denied. This does more than just assist in enhancing accountability but also addresses the issue of prejudice in lending practices among the populace. As Praxis AI pointed out, having SHAP generate individual explanations strengthens customer confidence and reduces litigation risk in financial institutions.

Oxley (2020) discussed how SHAP could be used to deconstruct black-box credit risk models and what characteristics drive risk assessment most. Unlike the classical metric of feature significance, SHAP additionally provides an institution with the ability to understand why a particular applicant was considered as either high risk or low risk, bridging the gap between the machine and the human perception of the applicant. Oxley explained that such transparency promotes well-intentioned applications of AI in finance, especially in cases where lack of interpretability can lead to discrimination or regulation violations.

Tandon et al. (2024) also studied the applicability of SHAP to the financial field and highlighted the importance of interpretability in helping to streamline operations. Their work described how, through the determination of the strongest factors in predicting loan defaults, banks can simplify their underwriting procedures and tailor their interventions to high-risk borrowers. Tandon and co-authors also noted that SHAP explanations tend to agree with the current economic theory, which supports the validity of AI models in the field of finance.

Arize AI (2024) introduced SHAP as one of the leading interpretability models to be used in fraud detection systems, especially credit card and insurance fraud. Models that are used in fraud detection are usually based on high-dimensional unbalanced data, making such data hard to interpret. SHAP copes with this issue by prioritizing features that are associated with suspicious behaviors, including the abnormal size of transactions or geolocation irregularities. As Arize AI noted, SHAP is not only more accurate in detecting fraud, but reportable audits, a compliance necessity under financial regulation, are also easier to achieve.

Contreras et al. (2024) proposed that the use of SHAP in finance may be further optimized by using group perturbation strategies like spectral zones. They reasoned that the correlation between groups of features (i.e., income, employment stability, credit history) yields more actionable information to regulators and financial analysts. Although these innovations are computationally demanding, they help to reduce redundancy and enhance interpretability to make decisions in finance.

SHAP helped identify behavioral drivers of adherence, such as appointment attendance or medication purchase frequency. These insights enabled healthcare providers to design behaviorally informed interventions, demonstrating SHAP's versatility across domains.

7. Challenges and Extensions of SHAP

Contreras et al. (2024) emphasized that one of the most significant challenges of SHAP is computational complexity, especially in deep neural networks and high-dimensional data. To overcome this, they suggested spectral zones-based SHAP, which clusters correlated features, which enhance scalability and interpretability without sacrificing accuracy. Likewise, Lundberg and Lee (2017) acknowledged that computing Shapley values in their exact form is NP-hard, and approximations such as KernelSHAP and TreeSHAP are partial in their relief but still consume a

substantial amount of computation. These disadvantages make SHAP less useful in real-time or scale.

Lundberg, Fu, and Lee (2020) highlighted the existence of an interpretability-usability trade-off, with SHAP being mathematically consistent, but the visual outputs might not be easily comprehensible to nonexpert users, turning the burden on models to stakeholders. The absence of standardized benchmarks with which SHAP can be assessed to measure its effectiveness was also reported in PMID: PMC10997568 (2024), demanding user-centered frameworks.

Practical issues are also applicable to enterprise environments. According to Opensense Labs (2025), SHAP cannot be truly helpful without visualization tools and decision-support systems. In NLP, Towards Data Science (2025) noted that fixed SHAP values do not represent contextual representations over time, which weakens their usefulness. Explanation overload was also documented in Muhammad et al. (2024) in genomic and imaging challenges.

Tandon et al. (2024) summarized that the future of SHAP depends on interdisciplinary cooperation in order to address computational limitations and improve human-friendly usability.

Discussion

The literature review confirms that SHAP (SHapley Additive exPlanations) is one of the most mathematically established and unified approaches in the explainable AI (XAI) area. The theoretical model created by Lundberg and Lee in 2017 to incorporate Shapley values into machine learning allowed model predictions to be broken down into additive, feature-specific contributions, which made SHAP unlike other interpretability methods that rely on heuristics. This solid foundation facilitates a certain degree of openness that creates confidence in high-stakes uses, particularly in contrast to substitute models such as LIME, which tend to be unstable and uncoherent globally.

Despite advantages, one of the central concerns brought up by the literature in all its incarnations is computer power: SHAP is computationally expensive, particularly when applied to deep learning networks, or high-dimensional inputs. Even though certain extensions such as spectral zones-based SHAP (Contreras et al., 2024) may be designed to seek efficiency optimization, scalability remains an issue with its use in a real-time or resource-constrained system. Where accuracy and speed are the most important factors like in finance and healthcare this limitation is very critical.

The second issue that is discovered is the lack of connection between the description of the technical model and the knowledge of what can be done by non-technical parties. SHAP provides mathematically sensible feature attributions, but even SHAP outputs must be translated with great care before being perceived as meaningful by practitioners, policymakers, or patients. That is why the collaboration of several fields, in other words, the use of AI skills and their application to the knowledge of the field to make the SHAP descriptions more accessible and useful is important.

Lastly, the literature indicates that despite its broad usage in the medical field and finance, other areas that need critical analysis, including autonomous systems, security in cyberspace, and social sciences, are underrepresented. Even standardization of benchmarks to assess the quality of an explanation in various application fields is lacking. Future studies need to focus on scalable estimates, inter-domain benchmarking, and domain-specific modifications to realize the full potential of SHAP as a general-purpose explainable AI tool.

Conclusion

In this review, the central role of SHAP is discussed due to the provision of a mathematically based and consistent framework on model interpretability. Shapley values Inspired by the ideas of cooperative game theory, SHAP offers additive attributions of features which makes it both local and globally consistent, a trait it shares with additive attributions of other heuristic-based methods such as LIME. SHAP has become one of the most popular frameworks to shed light on the black box of AI because of its sound theoretical basis and broad applications.

Literature shows the importance of SHAP in high stakes areas like healthcare, finance, and customer analytics where honesty, equity, and responsibility are vital. SHAP has enhanced transparency in the diagnostics, as well as regulatory compliance and decision-making, by providing a detailed feature-level insight. Its application in both natural language processing and marketing further underscores its ability to improve both the reliability of the model and human comprehension. However, challenges remain. Scalability is still limited by computational inefficiency, especially in deep learning and high-dimensional data. The disparity between technical production and understandable interpretation also highlights that better communication of SHAP findings to non-technical stakeholders is required.

These constraints notwithstanding, the further development of SHAP, by extensions, hybrid techniques, and domain-specific specializations, provides opportunities to address existing obstacles. The future direction of research should aim at building effective approximations, introducing standard benchmarks, and extending SHAP to new fields, including cybersecurity, autonomous systems and social sciences. Finally, SHAP is not only a potent interpretability framework but also a starting point toward building trust, accountability, and transparency in AI systems in various fields.

References:

1. Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765–4774. <https://doi.org/10.48550/arXiv.1705.07874>
2. Muhammad, D., et al. (2024). Unveiling the black box: A systematic review of explainable AI methods in medical image analysis. *Frontiers in Artificial Intelligence*, 7, Article 11382209.
3. Contreras, J., et al. (2024). Spectral zones-based SHAP/LIME: Enhancing interpretability of deep learning models with group perturbations. *Analytical Chemistry*, 96(18), 7329-7337. <https://doi.org/10.1021/acs.analchem.4c02329>
4. Phalkey, M. (2023, April 17). Model explainability and interpretability using SHAP. *almaBetter*.
5. Lundberg, S. M., Fu, H., & Lee, S.-I. (2020). A perspective on explainable artificial intelligence methods: SHAP and LIME. *Information Fusion*, 58, 82-86. <https://arxiv.org/abs/2305.02012>
6. Opensense Labs. (2025, January 30). Explainable AI tools: SHAP's power in AI. <https://opensenselabs.com/blog/explainable-ai-tools>
7. Arize AI. (2024, January 10). Explainability in machine learning: Top techniques including SHAP. <https://arize.com/blog-course/explainability-techniques-shap>
8. Tandon, A., Ansari, R., Tandon, C., Shah, R., & Saranya, A. (2024). Illuminating the black box: Explainable AI (XAI). *International Journal of Advanced Computer Engineering and Networking*, 10(3), 45-56.

9. PMID: PMC10997568. (2024). Interpreting artificial intelligence models: A systematic review on SHAP explainability techniques. *Journal XYZ*, 12(4), 123-134.
10. Praxis AI. (2025, January 19). Illuminating the 'black box' with explainable AI (XAI). <https://praxis.ac.in/illuminating-the-black-box-with-explainable-ai-xai/>
11. International Journal of Science and Research Archive. (2025, April 20). Explainable AI for enhanced customer behavior prediction and trust. *IJSRA*, 15(1), 1-14.
12. Oxley, M. (2020). Machine learning model explanation: Illuminating the black box. University of Edinburgh Business School.
13. Lundberg, S. M., & Lee, S.-I. (2017). SHAP: SHapley additive explanations. *Advances in Neural Information Processing Systems*. [PDF document]
14. Towards Data Science. (2025, January 14). Illuminating the black box of textual GenAI using SHAP. <https://towardsdatascience.com/illuminating-the-black-box-of-ai-ddea07e65c35/>
15. Muhammad, D., et al. (2024). Systematic review on SHAP and other explanation methods for AI in healthcare. *Journal of Healthcare Informatics Research*, 8(2), 345-359.