

On Generating Monolithic and Model Reconciling Explanations in Probabilistic Scenarios (Abstract Reprint)

Stylios Loukas Vasileiou¹, William Yeoh², Alessandro Previti³, Tran Cao Son¹

¹New Mexico State University, United States

²Washington University in St. Louis, United States

³Ericsson Research, Sweden

Abstract Reprint. This is an abstract reprint of the journal article by Vasileiou, Yeoh, Previti, and Son (2025).

Abstract

Explanation generation frameworks aim to make AI systems decisions transparent and understandable to human users. However, generating explanations in uncertain environments characterized by incomplete information and probabilistic models remains a significant challenge. In this paper, we propose a novel framework for generating probabilistic monolithic explanations and model reconciling explanations. Monolithic explanations provide self-contained reasons for an explanandum without considering the agent receiving the explanation, while model reconciling explanations account for the knowledge of the agent receiving the explanation. For monolithic explanations, our approach integrates uncertainty by utilizing probabilistic logic to increase the probability of the explanandum. For model reconciling explanations, we propose a framework that extends the logic-based variant of the model reconciliation problem to account for probabilistic human models, where the goal is to find explanations that increase the probability of the explanandum while minimizing conflicts between the explanation and the probabilistic human model. We introduce explanatory gain and explanatory power as quantitative metrics to assess the quality of these explanations. Further, we present algorithms that exploit the duality between minimal correction sets and minimal unsatisfiable sets to efficiently compute both types of explanations in probabilistic contexts. Extensive experimental evaluations on various benchmarks demonstrate the effectiveness and scalability of our approach in generating explanations under uncertainty.

References

Vasileiou, S. L.; Yeoh, W.; Previti, A.; and Son, T. C. 2025. On Generating Monolithic and Model Reconciling Explanations in Probabilistic Scenarios. *Journal of Artificial Intelligence Research*, 84: 5:1–5:40.