

CONFIDERA: CONFormal Interpretable-by-Design score function for Explainable and Reliable Artificial Intelligence

*Original*

CONFIDERA: CONFormal Interpretable-by-Design score function for Explainable and Reliable Artificial Intelligence / Narteni, S., Carlevaro, A., Muselli, M., Dabbene, F., Mongelli, M.. - ELETTRONICO. - 204:(2023), pp. 485-487. (The 12th Symposium on Conformal and Probabilistic Prediction with Applications (COPA 2023) Limassol (CY) 13-15 September 2023).

*Availability:*

This version is available at: 11583/2982243 since: 2023-09-22T07:04:03Z

*Publisher:*

Proceedings of Machine Learning Research

*Published*

DOI:

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)

# CONFIDERAi: CONFormal Interpretable-by-Design score function for Explainable and Reliable Artificial Intelligence

Sara Narteni<sup>1,†</sup>

Alberto Carlevaro<sup>1,†</sup>

Fabrizio Dabbene<sup>1</sup>

Marco Muselli<sup>1,2</sup>

Maurizio Mongelli<sup>1</sup>

SARA.NARTENI@IEIIT.CNR.IT

ALBERTO.CARLEVARO@IEIIT.CNR.IT

FABRIZIO.DABBENE@IEIIT.CNR.IT

MARCO.MUSELLI@IEIIT.CNR.IT

MAURIZIO.MONGELLI@IEIIT.CNR.IT

<sup>1</sup> *CNR-IEIIT, 10129, Turin, Italy*

<sup>2</sup> *Rulex Innovation Labs, Rulex Inc., 16122 Genoa, Italy*

† *S. Narteni and A. Carlevaro contributed equally to the development of the article. (Corresponding authors: S.Narteni, A. Carlevaro.)*

**Editor:** Harris Papadopoulos, Khuong An Nguyen, Henrik Boström and Lars Carlsson

## Abstract

The concept of trustworthiness has been declined in different ways in the field of artificial intelligence, but all its definitions agree on two main pillars: explainability and conformity. In this extended abstract, our aim is to give an idea on how to merge these concepts, by defining a new framework for conformal rule-based predictions. In particular, we introduce a new score function for rule-based models, that leverages on rule relevance and geometrical position of points from rule classification boundaries.

**Keywords:** XAI, conformal safety sets, novel score function, conformal prediction.

## 1. Introduction

Literature around combination of eXplainable AI (XAI) and conformal prediction (CP) has recently gained popularity but it still remains little investigated in research. Some relevant approaches proposed so far investigated conformal prediction for XAI models (Bhattacharyya, 2011; Johansson et al., 2014, 2018, 2022), but to the best of our knowledge, no previous study of this type addressed score functions and quantile, tailored for rule-based models.

For this reason, we propose CONFIDERAi, an innovative approach, based on a new score function, to build conformal prediction of rule-based models. The rationale behind the approach is the combination of the global properties of decision rules (i.e., their covering and error) and the geometrical position of the points inside rule boundaries. The resulting prediction set leads to a restricted *conformal safety set*, i.e., the set of points for which the underlying XAI model performs with probabilistic guarantees.

## 2. CONFIDERAi

**Conformal Safety Set.** CSS allows to insert CP in a more safety-based context. For any input feature  $\mathbf{x} \in \mathcal{X}$  and any label  $y \in \mathcal{Y}$ , given a *prediction set* at *level of confidence*  $1 - \varepsilon$ ,  $\varepsilon \in (0, 1)$ ,

$$\mathcal{C}(\mathbf{x}) = \{y \mid s(\mathbf{x}, y) \leq s_\varepsilon\} \in 2^{\mathcal{Y}}, \quad (1)$$

where  $s_\varepsilon$  is the  $1 - \varepsilon$  quantile of the score values computed on a *calibration set*, CSS is defined as a subset of the input feature space in which probabilistic safety guarantees can be provided to the machine learning (ML) model:

$$\mathcal{S}_\varepsilon = \{\mathbf{x} \mid \Pr \{y \in \mathcal{C}(\mathbf{x})\} \geq 1 - \varepsilon, \forall y \in \mathcal{Y}\} = \{\mathbf{x} \mid s(\mathbf{x}, y) \leq s_\varepsilon, \forall y \in \mathcal{Y}\}. \quad (2)$$

**Novel Score Function.** An innovative score function suitable to find  $\mathcal{S}_\varepsilon$  for rule-based models is designed as follows:

$$s(\mathbf{x}, y) \doteq \sum_{r_k \in \mathcal{R}_\mathbf{x}^y} \tau(\mathbf{x}, r_k)(1 - R(r_k)) \quad \text{where} \quad \tau(\mathbf{x}, r_k) = \frac{1}{1 + e^{-1/\gamma}}, \quad (3)$$

with  $\gamma = \gamma(\mathbf{x}, r_k)$  the minimum of the Euclidean distances between the point  $\mathbf{x}$  covered by the rule  $r_k$  and each side of the rule boundary. Values of  $\tau$  closer to 1 thus encode higher proximity to rule boundary and probability of misclassification. Rule relevance  $R(r_k)$  is also accounted Ferrari et al. (2022). The sum is over rules belonging to the set  $\mathcal{R}_\mathbf{x}^y$  of rules covered by  $\mathbf{x}$  predicting class  $y$ .

**Preliminary Results.** A further classifier is trained to distinguish *conformal* and *non conformal* points, i.e., individuate a *Conformal Safety Region*.

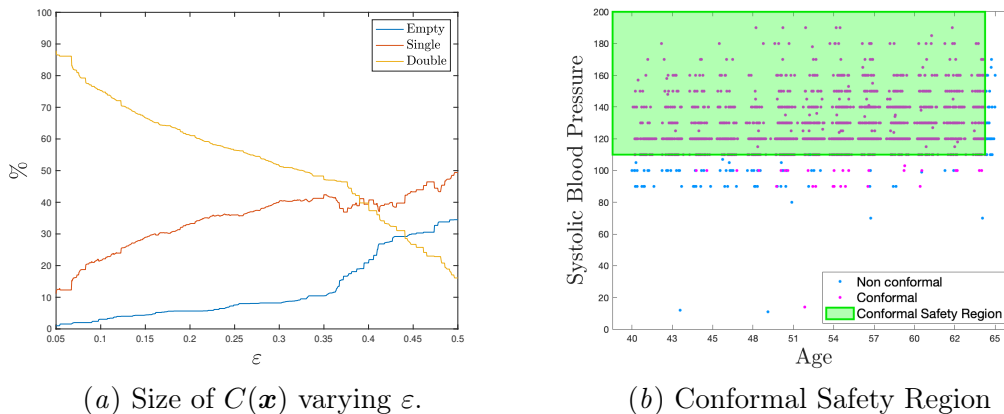


Figure 1: Risk Factors for Cardiovascular Heart Disease (CHD) dataset <sup>\*</sup>.

## Acknowledgements

This work was supported in part by REXASI-PRO H-EU project, call HORIZON-CL4-2021-HUMAN-01-01, Grant agreement ID: 101070028. The work was also supported by Future Artificial Intelligence Research (FAIR) project, Recovery and Resilience Plan ("Piano Nazionale di Ripresa e Resilienza"), Spoke 3 - Resilient AI.

<sup>\*</sup><https://www.kaggle.com/datasets/thedevastator/exploring-risk-factors-for-cardiovascular-diseas>

## References

- Siddhartha Bhattacharyya. Confidence in predictions from random tree ensembles. In *2011 IEEE 11th International Conference on Data Mining*, pages 71–80, 2011. doi: 10.1109/ICDM.2011.41.
- Enrico Ferrari, Damiano Verda, Nicolò Pinna, and Marco Muselli. A novel rule-based modeling and control approach for the optimization of complex water distribution networks. In *Advances in System-Integrated Intelligence: Proceedings of the 6th International Conference on System-Integrated Intelligence (SysInt 2022), September 7-9, 2022, Genova, Italy*, pages 33–42. Springer, 2022.
- Ulf Johansson, Rikard König, Henrik Linusson, Tuve Löfström, and Henrik Boström. Rule extraction with guaranteed fidelity. In *Artificial Intelligence Applications and Innovations: AIAI 2014 Workshops: CoPA, MHDW, IIVC, and MT4BD, Rhodes, Greece, September 19-21, 2014. Proceedings 10*, pages 281–290. Springer, 2014.
- Ulf Johansson, Henrik Linusson, Tuve Löfström, and Henrik Boström. Interpretable regression trees using conformal prediction. *Expert systems with applications*, 97:394–404, 2018.
- Ulf Johansson, Cecilia Sönströd, Tuve Löfström, and Henrik Boström. Rule extraction with guarantees from regression models. *Pattern Recognition*, 126:108554, 2022.