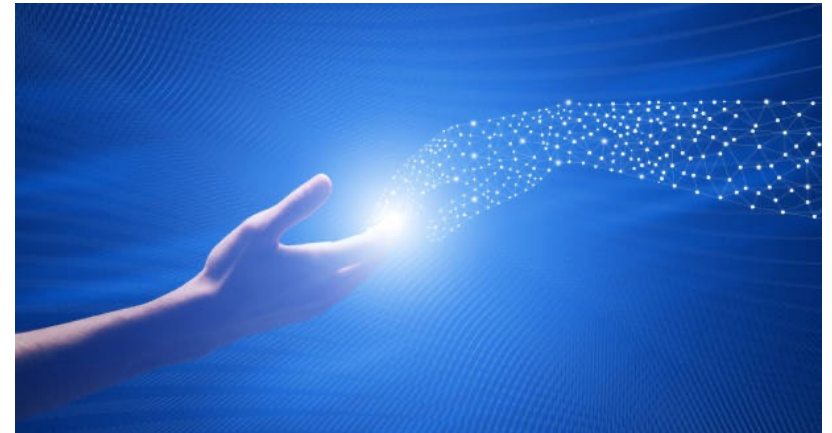


**“CONFIDERA1: CONFormal Interpretable-by-Design score function for
Explainable and Reliable Artificial Intelligence”**

Sara Narteni, Alberto Carlevaro, Fabrizio Dabbene, Marco Muselli, Maurizio Mongelli
Cnr-Istituto di Elettronica e di Ingegneria dell'Informazione e delle Telecomunicazioni
sara.narteni@ieiit.cnr.it, alberto.carlevaro@ieiit.cnr.it

OUTLINE

- Motivation and background
- CONFIDERAIA
 - Conformal Safety Set
 - Novel score function for rule-based models
- Preliminary applications
- Conclusions and future steps



MOTIVATION AND BACKGROUND



- Increased diffusion of **trustworthy AI**
 1. Need of *transparency*: **eXplainable AI (XAI)**
 2. Need of *safety guarantees*
- Strong **conformal prediction** frameworks
 1. (non)-conformity measure + p-values [1]
 2. score functions + quantile [2]
- Little investigation on CP for XAI models (with the 1st technique)

Our work: combining XAI and CP
by designing a new score function
for rule-based models

Picture credit: European Commission, Directorate-General for Communications Networks, Content and Technology, *Ethics guidelines for trustworthy AI*, Publications Office, 2019, <https://data.europa.eu/doi/10.2759/346720>

[1] [Glenn Shafer, Vladimir Vovk](https://arxiv.org/abs/0706.3188), a tutorial on conformal prediction, arXiv, 2007, <https://arxiv.org/abs/0706.3188>

[2] [Anastasios N. Angelopoulos, Stephen Bates](https://arxiv.org/abs/2107.07511), A Gentle Introduction to Conformal Prediction and Distribution-Free Uncertainty Quantification, arXiv, 2021, <https://arxiv.org/abs/2107.07511>

FROM CONFORMAL PREDICTION TO CONFORMAL SAFETY SETS

Machine Learning model

$$\hat{f}(x)_y : \mathcal{X} \rightarrow \mathcal{Y} \sim$$

Score function

$$s(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$

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

Size of *calibration set*

$$\frac{\lceil (n_c + 1)(1 - \varepsilon) \rceil}{n_c} \text{ quantile}$$

$$\mathcal{S}_\varepsilon = \left\{ x \mid s(x, y) \leq s_\varepsilon, \forall y \in \mathcal{Y} \right\}$$

Conformal Safety Set
(all labels)

$$x_i \rightarrow \begin{cases} +1 & \text{if } x_i \in \mathcal{S}_\varepsilon \\ -1 & \text{otherwise} \end{cases}$$

Binary
Classifier

Conformal Safety Region

New dataset labelling

FROM CONFORMAL PREDICTION TO CONFORMAL SAFETY SETS

Machine Learning model

$$\hat{f}(\mathbf{x})_y : \mathcal{X} \rightarrow \mathcal{Y} \sim$$

Score function

$$s(\mathbf{x}, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$

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

Size of *calibration set*

$$\frac{\lceil (n_c + 1)(1 - \varepsilon) \rceil}{n_c} \text{ quantile}$$

$$\mathcal{S}_\varepsilon = \left\{ \mathbf{x} \mid s(\mathbf{x}, +1) \leq s_\varepsilon, s(\mathbf{x}, 0) > s_\varepsilon \right\}$$

Ongoing extension:
focus on a **target class** (i.e., +1)

$$\mathbf{x}_i \rightarrow \begin{cases} +1 & \text{if } \mathbf{x}_i \in \mathcal{S}_\varepsilon \\ -1 & \text{otherwise} \end{cases}$$

Binary
Classifier

Conformal Safety Region

New dataset labelling

NOTATION FOR RULE-BASED MODELS

$$\mathcal{T} = \{(\mathbf{x}_j, y_j)\}_{j=1}^N, \text{ with } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d \text{ and } y \in \{0, 1\}$$

$$\mathcal{R} = \{r_k\}_{k=1}^M \longrightarrow \text{if premise then consequence}$$

The *premise* constitutes the antecedent of the rule and is a logical conjunction (\wedge) of conditions c_{i_k} , with $i_k = 1_k, \dots, N_k$

$$\left\{ \begin{array}{l} 1. x_{\pi(i)} \geq l_{i_k} \\ 2. x_{\pi(i)} \leq u_{i_k} \\ 3. l_{i_k} \leq x_{\pi(i)} \leq u_{i_k} \end{array} \right.$$

$$\left. \begin{array}{l} C(r_k) = \frac{TP(r_k)}{TP(r_k) + FN(r_k)} \\ E(r_k) = \frac{FP(r_k)}{TN(r_k) + FP(r_k)} \end{array} \right\} R(r_k) = C(r_k) \cdot (1 - E(r_k)) \text{ Rule relevance}$$

NOVEL SCORE FUNCTION FOR RULE-BASED MODELS

$\mathbf{x} = (x_1, x_2)$: 2D input point

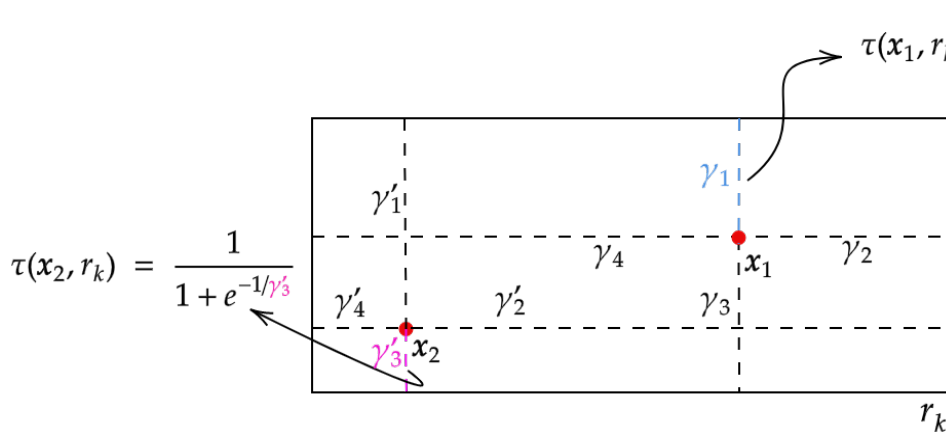
y : output label

r_k : k_{th} rule of the model

$\gamma = \gamma(\mathbf{x}, r_k)$: minimum of Euclidean distances between \mathbf{x} and rule boundary sides

$$s(\mathbf{x}, y) \doteq \sum_{r_k \in \mathcal{R}_x^y} \tau(\mathbf{x}, r_k) \underbrace{(1 - R(r_k))}_{\text{rule relevance}}$$

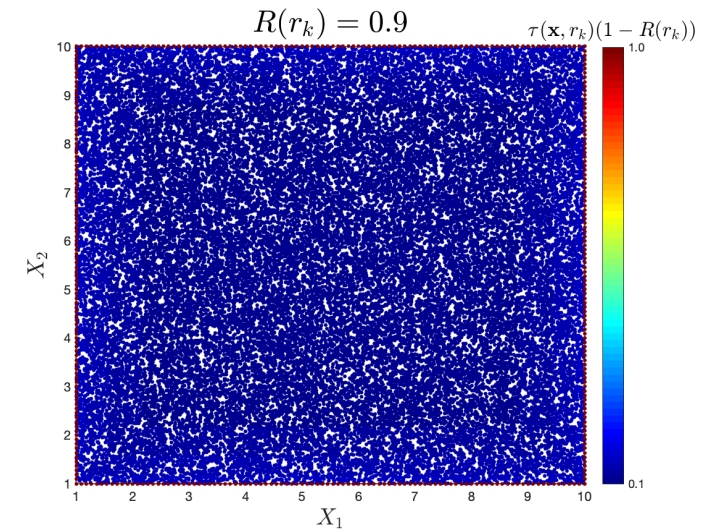
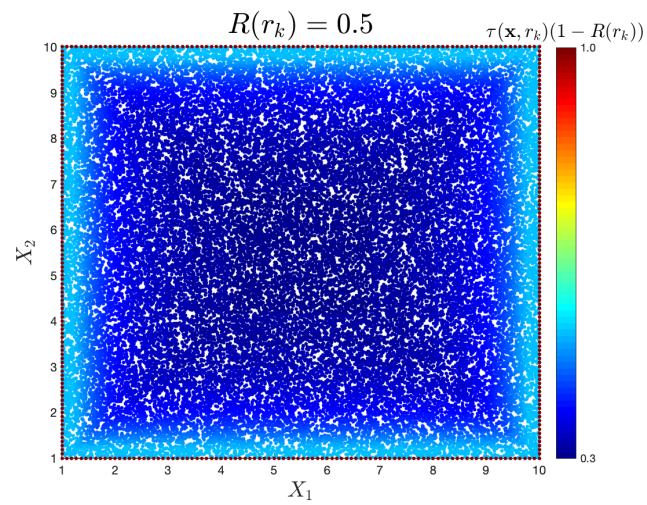
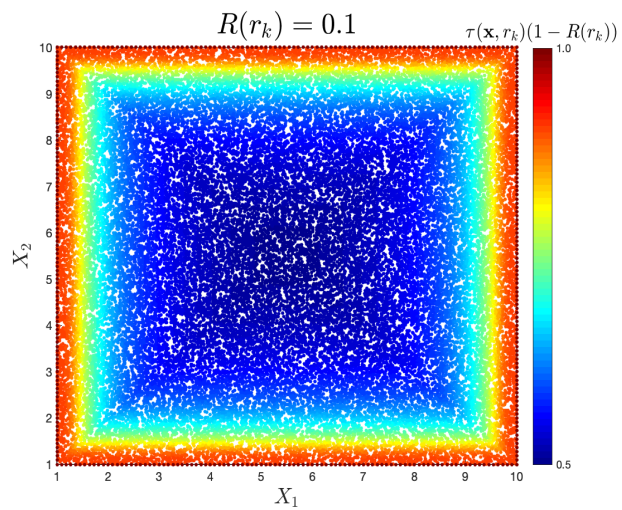
rule predicting y and verified by \mathbf{x} position of point within rule boundary



$$\tau(\mathbf{x}, r_k) = \frac{1}{1 + e^{-1/\gamma}}$$

$$\tau(\mathbf{x}_2, r_k) > \tau(\mathbf{x}_1, r_k)$$

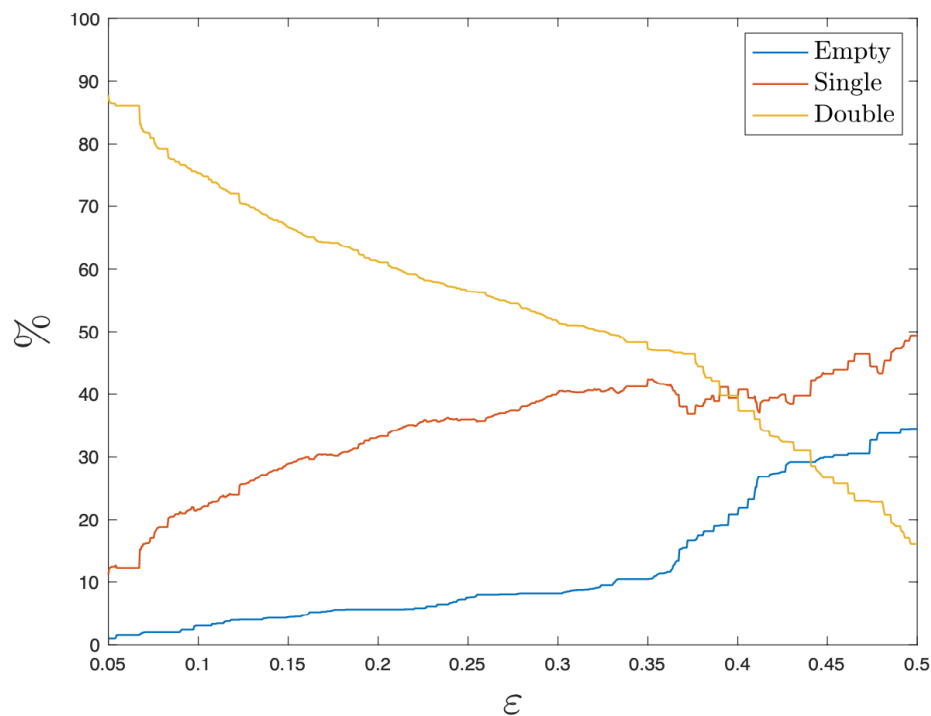
ILLUSTRATIVE EXAMPLE



rule relevance increase
constant $\tau(\mathbf{x}, r_k)$

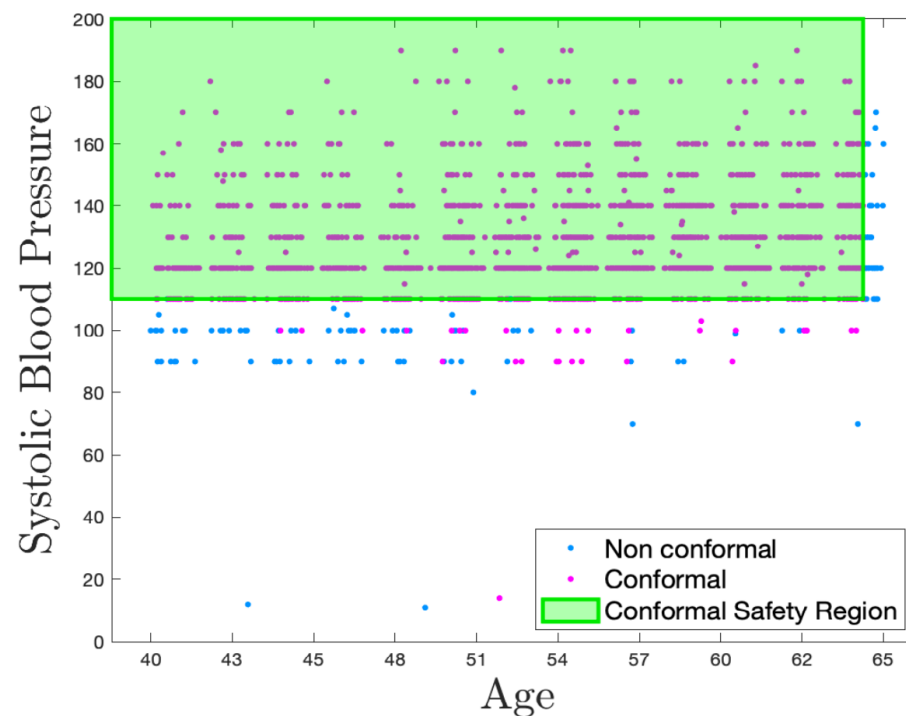
PRELIMINARY RESULTS

Efficiency:
average size of the prediction region



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

Conformal Safety Region



CONCLUSIONS

- New **score function** to perform CP on rule-based (transparent-by-design) XAI models in 2D, accounting for:
 - Performance ability
 - Geometry
- Step beyond pure CP:
 - definition of **conformal safety sets (CSS) and regions**
- Future (ongoing) works:
 - > 2D extension
 - Refinement of CSS definition towards safety preservation



THANK YOU!
Q&A Time