

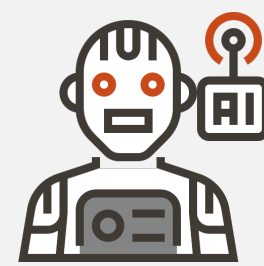
GLocalX - From Local to Global Explanations of Black Box AI Models

Mattia Setzu¹ Riccardo Guidotti¹ Anna Monreale¹ Franco Turini¹ Dino Pedreschi¹ Fosca Giannotti²

¹University of Pisa, Italy ²ISTI-CNR, Italy

The Local to Global explanation problem

Local explanations tend to be faithful yet complex, while global explanations tend to be more general yet less faithful. **Can we retain the faithfulness of local explanations and the simplicity of global ones?** In a *Local to Global* setting [2] we are given:



black-box model



set of local explanations



handful of data

and aim to infer a set of global explanations from the local ones.

The GLocalX algorithm

GLocalX takes after aggregating hierarchical clustering algorithms. It yields global explanations by iteratively

1. identifying similar sets of explanations,
2. merging them,
3. validating their merge.

The resulting dendrogram is a binary tree whose nodes hold local explanations in its leaves, and increasingly global explanations as the nodes near the root. We extract a set of explanations of size α from the root node(s).

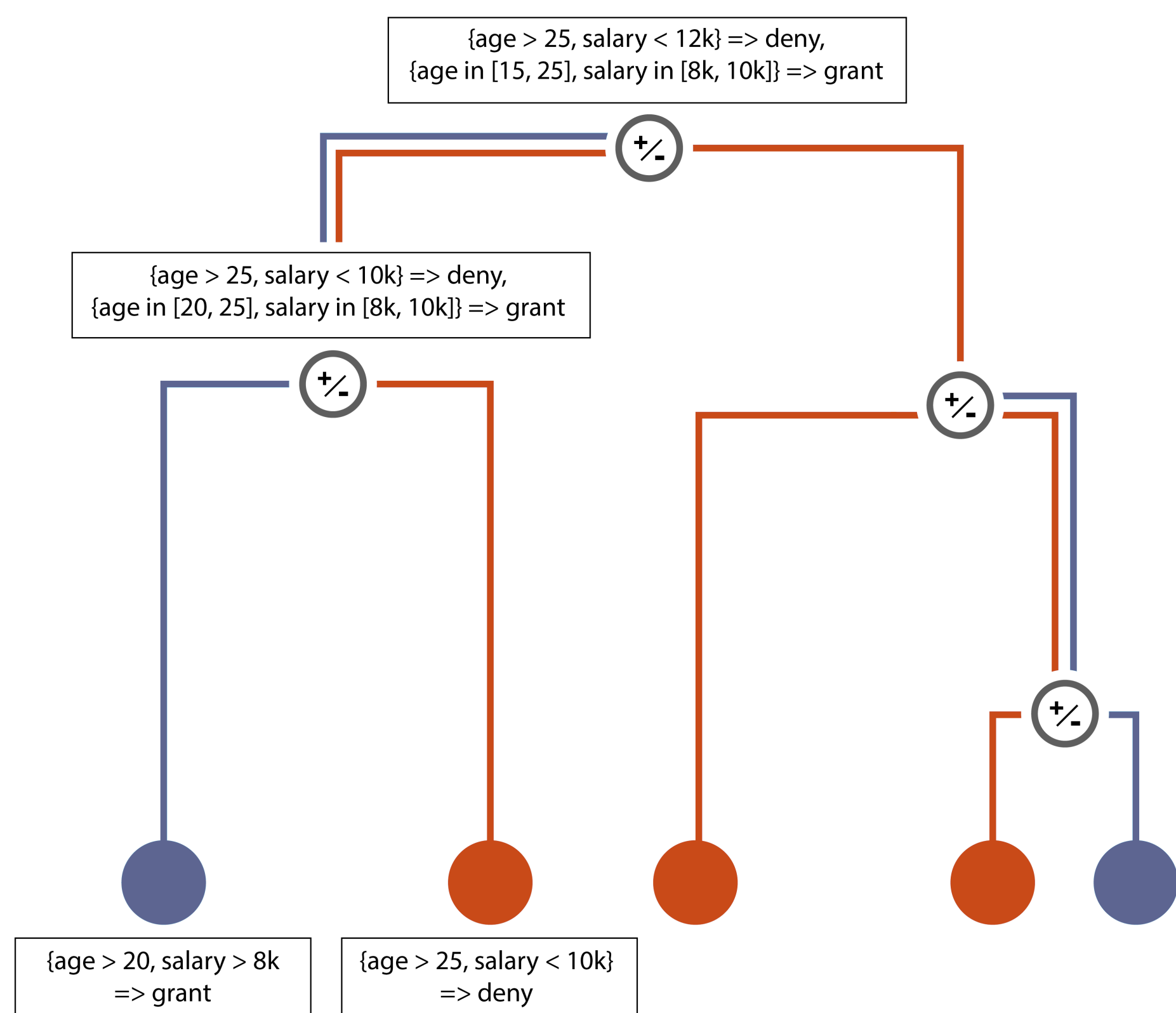


Figure 1. An example of merge in the dendrogram.

Here we focus on *axis-parallel decision rules* as local explanations, locally inferred by LORE [1].

Finding Similar Sets of Explanations

Similarity is measured as Jaccard similarity of explanation *coverage*, that is the larger the ratio of records explained by two explanations, the higher their similarity and viceversa.

$$\text{sim}(E_1, E_2) = \frac{\text{cov}(E_1, X) \cap \text{cov}(E_2, X)}{\text{cov}(E_1, X) \cup \text{cov}(E_2, X)}$$

When dealing with sets of explanations, the coverage of a set of explanations E can be:

linkage	explanation(s) selected	coverage
min	explanation e with minimum coverage	$\text{cov}(e, X), e : \arg \min_{e \in E} \text{cov}(e, X)$
max	explanation e with maximum coverage	$\text{cov}(e, X), e : \arg \max_{e \in E} \text{cov}(e, X)$
full	the union of all explanations	$\bigcup_{e \in E} \text{cov}(e, X)$

Merging explanations

An explanation is a polyhedra in the feature space, each premise $P_i : [a, b]$ identifying a continuous interval $[a, b]$ in the feature space of feature i .

$$\underbrace{\text{age} \geq 20}_{P_{\text{age}=[20, +\infty)}} , \underbrace{\text{salary} \geq 8k}_{P_{\text{salary}=[8k, +\infty)}} \rightarrow \text{grant}$$

Merge is bipartite:

- the join operator (\oplus) generalizes concordant explanations, increasing their coverage and possibly decreasing their fidelity;
- the cut operator (\ominus) localizes discordant explanations, decreasing their coverage, possibly increasing their fidelity.

	Case	Merging	Merged
\oplus	join		
\oplus	join empty		
\ominus	cut left		
\ominus	cut right		
\ominus	cut in between		
\ominus	cut everything		

Table 1. Visual representation of the join (\oplus) and cut (\ominus) operators.

Accepting and Rejecting Merges

Not all merges are created equal. A merge is accepted if it yields a faithful and simple explanations [3]:

$$\text{BIC}(E) = \text{fid}(E, X, \tilde{Y}) + \frac{1}{|E|} \sum_{e \in E} \text{len}(e) + \frac{|E|}{n}$$

Local to Global vs Natively Global and Natively Local

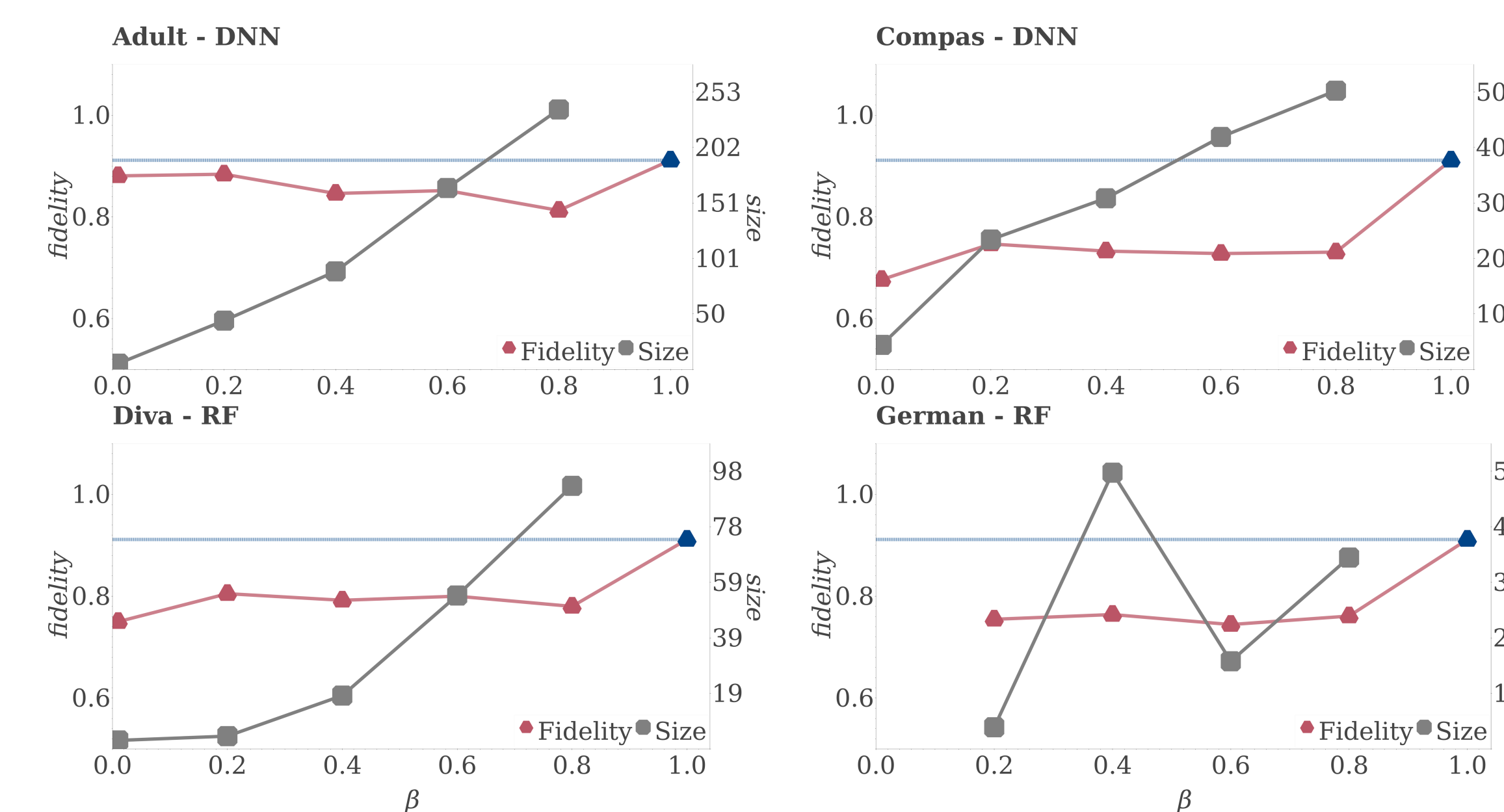
Table 2. GLocalX (trained on synthetic data) vs natively global models and a baseline of local explanations union.

Deep Network	Adult, $\alpha = 10$			Compas, $\alpha = 5$		
	Fidelity	Size	Length	Fidelity	Size	Length
GLocalX	.912	<u>5</u>	6.00 ± 1.0	.759	<u>3</u>	2.33 ± 0.4
GLocalX*	.880	<u>10</u>	6.30 ± 2.5	.756	<u>6</u>	4.50 ± 0.9
CPAR	<u>.929</u>	100	3.78 ± 2.4	<u>.821</u>	69	3.11 ± 1.4
Decision Tree	.917	1068	7.22 ± 1.9	.789	1014	6.00 ± 1.8
Pruned Decision Tree	.908	28	2.71 ± 0.8	.780	30	2.33 ± 0.6
Union	.880	6838	3.86 ± 2.2	.627	1515	4.65 ± 1.7
Random Forest		Diva, $\alpha = 25$		German, $\alpha = 2$		
GLocalX	<u>.854</u>	<u>26</u>	3.26 ± 0.9	.786	<u>2</u>	3.00 ± 1.0
GLocalX*	.848	<u>26</u>	3.88 ± 1.3	.766	<u>2</u>	5.87 ± 2.4
CPAR	.850	221	<u>2.03 ± 1.0</u>	.773	18	<u>2.33 ± 1.4</u>
Decision Tree	.853	976	10.63 ± 3.9	<u>.830</u>	76	4.50 ± 1.7
Pruned Decision Tree	.836	28	3.21 ± 0.9	.796	28	2.78 ± 0.8
Union	.794	2013	2.90 ± 1.0	.766	210	5.62 ± 2.4

What if I have less rules available?

Extracting local explanations can be time-consuming, can we reduce the input size?

Figure 2. GLocalX performance with $\beta \in 1, 20, 40, 60, 80$ percentage of input rules, randomized in 10 trials, and $\alpha_q = 50$.



References

- [1] R. Guidotti, A. Monreale, F. Giannotti, D. Pedreschi, S. Ruggieri, and F. Turini. Factual and counterfactual explanations for black box decision making. *IEEE Intelligent Systems*, 34(6):14–23, 2019.
- [2] D. Pedreschi, F. Giannotti, R. Guidotti, A. Monreale, S. Ruggieri, and F. Turini. Meaningful explanations of black box ai decision systems. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 9780–9784, 2019.
- [3] E. Wit, E. v. d. Heuvel, and J.-W. Romeijn. ‘all models are wrong...’: an introduction to model uncertainty. *Statistica Neerlandica*, 66(3):217–236, 2012.