

Enhancing Interpretability of Black Box Models by means of Local Rules

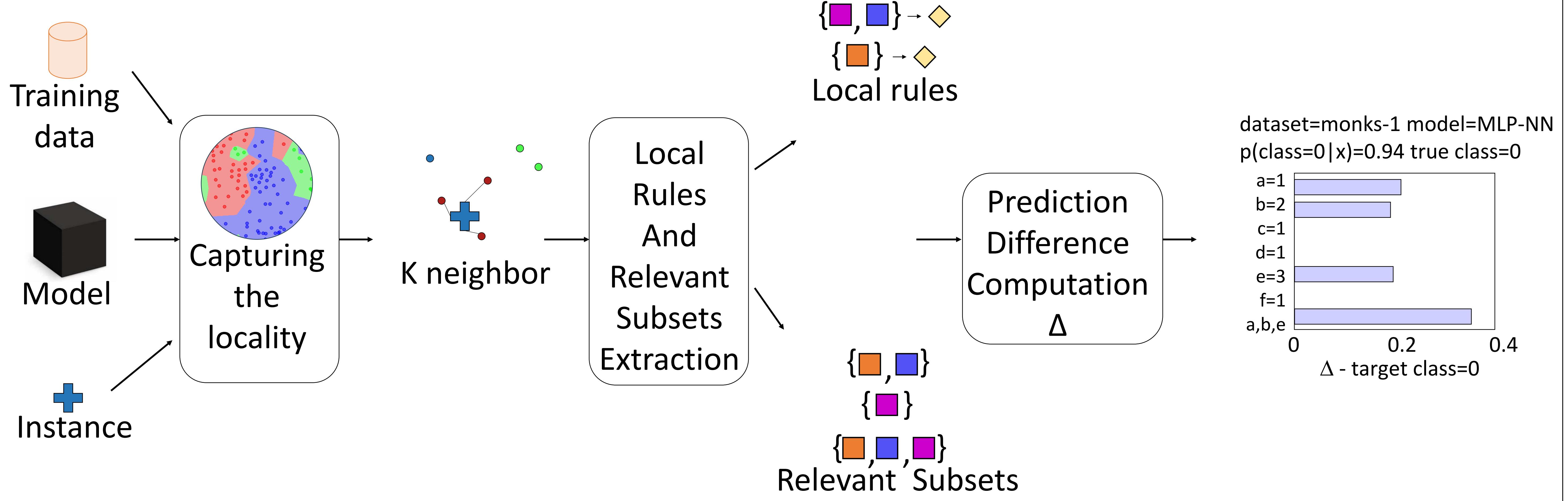
eliana.pastor@polito.it

Eliana Pastor and Elena Baralis

Politecnico di Torino, Italy



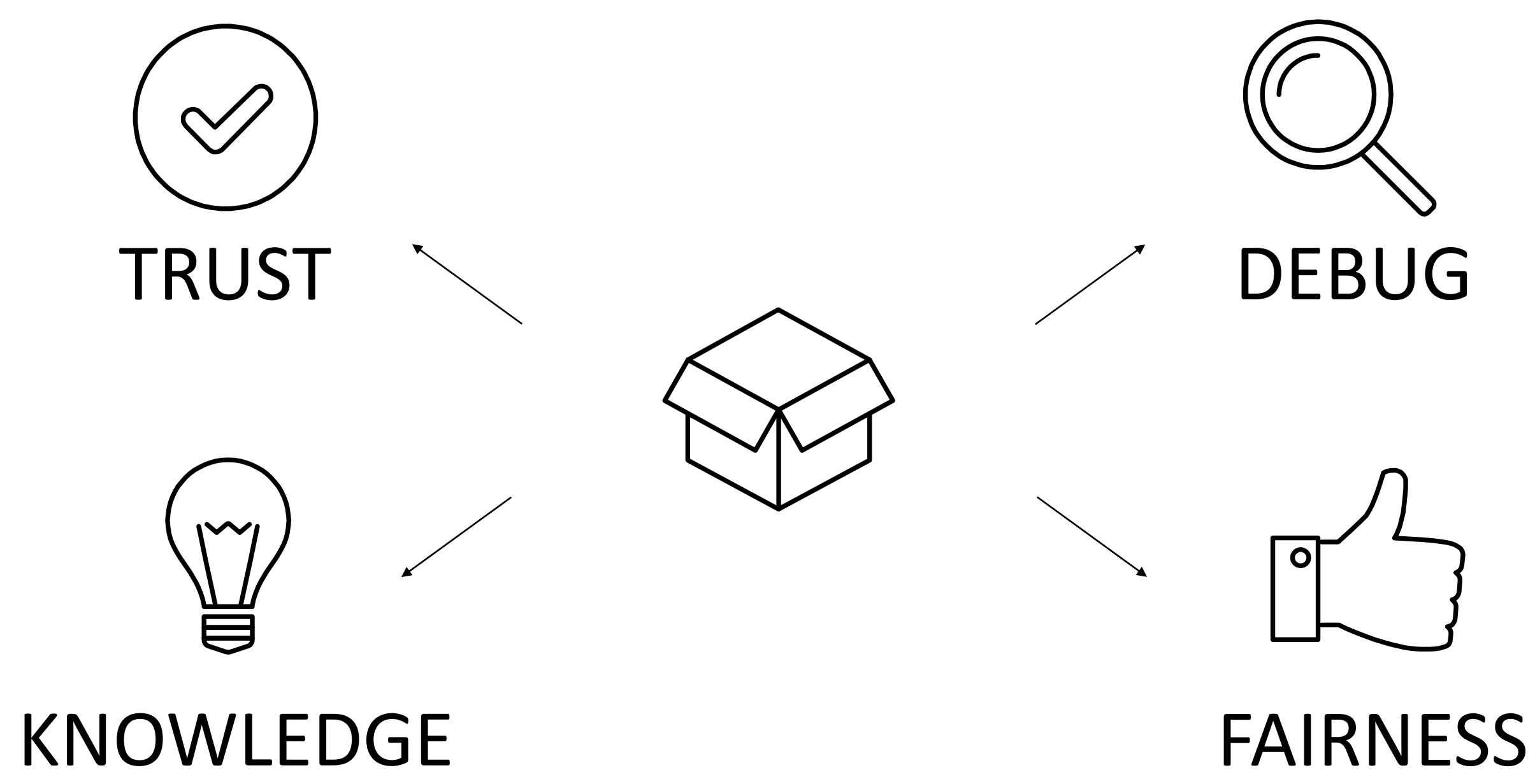
LACE explanation method



Enhancing interpretability

Black box models → hide from users the reasons behind predictions

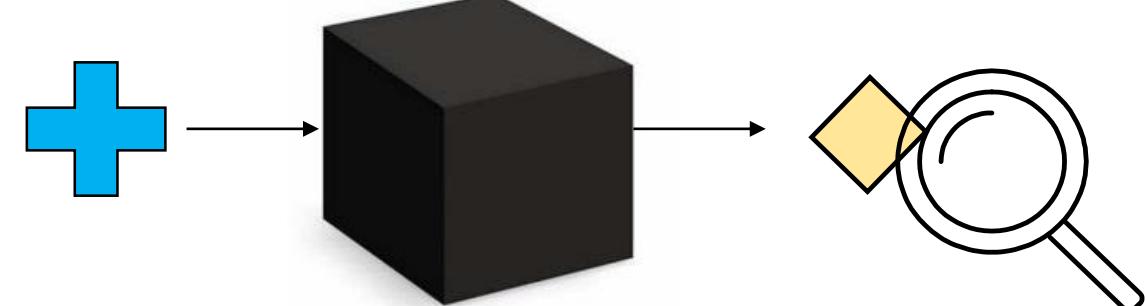
Why explainable AI?



Demand from institutions → GDPR

- “meaningful information about the logic involved”
- “to ensure fair and transparent processing”

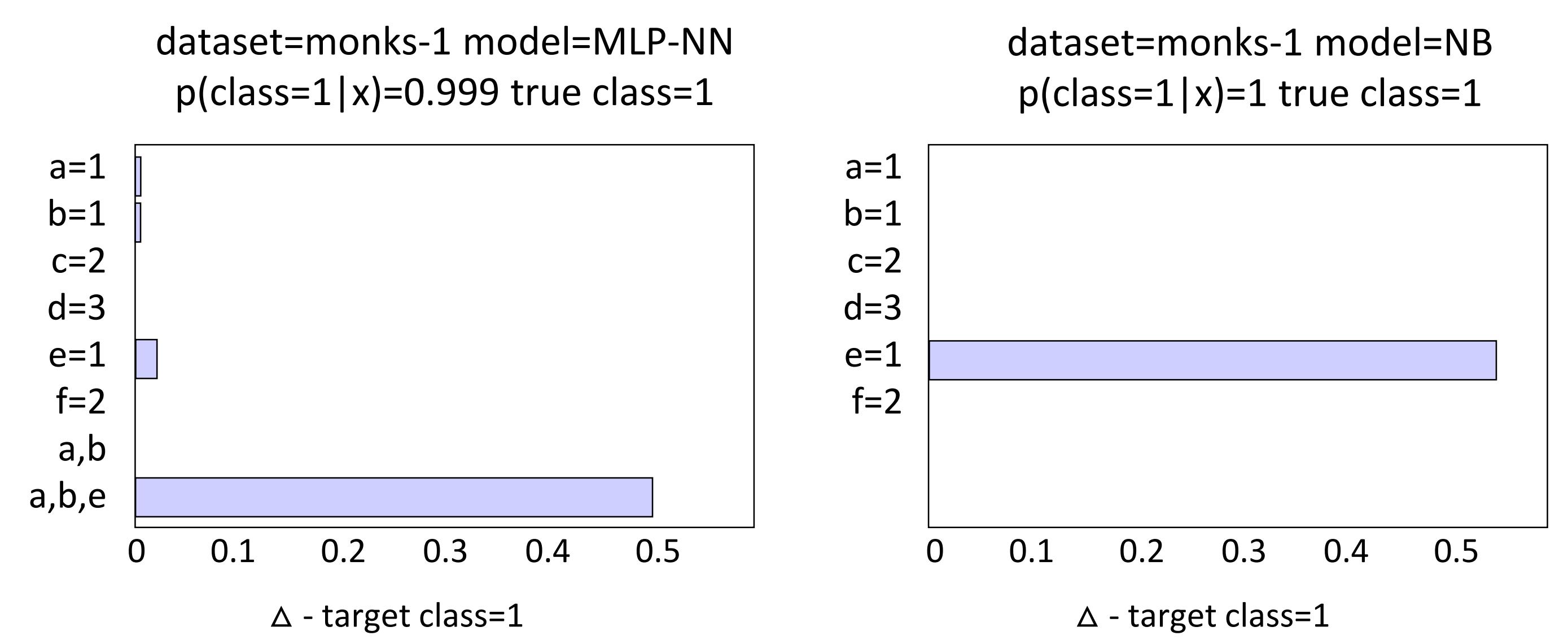
LACE Explanations



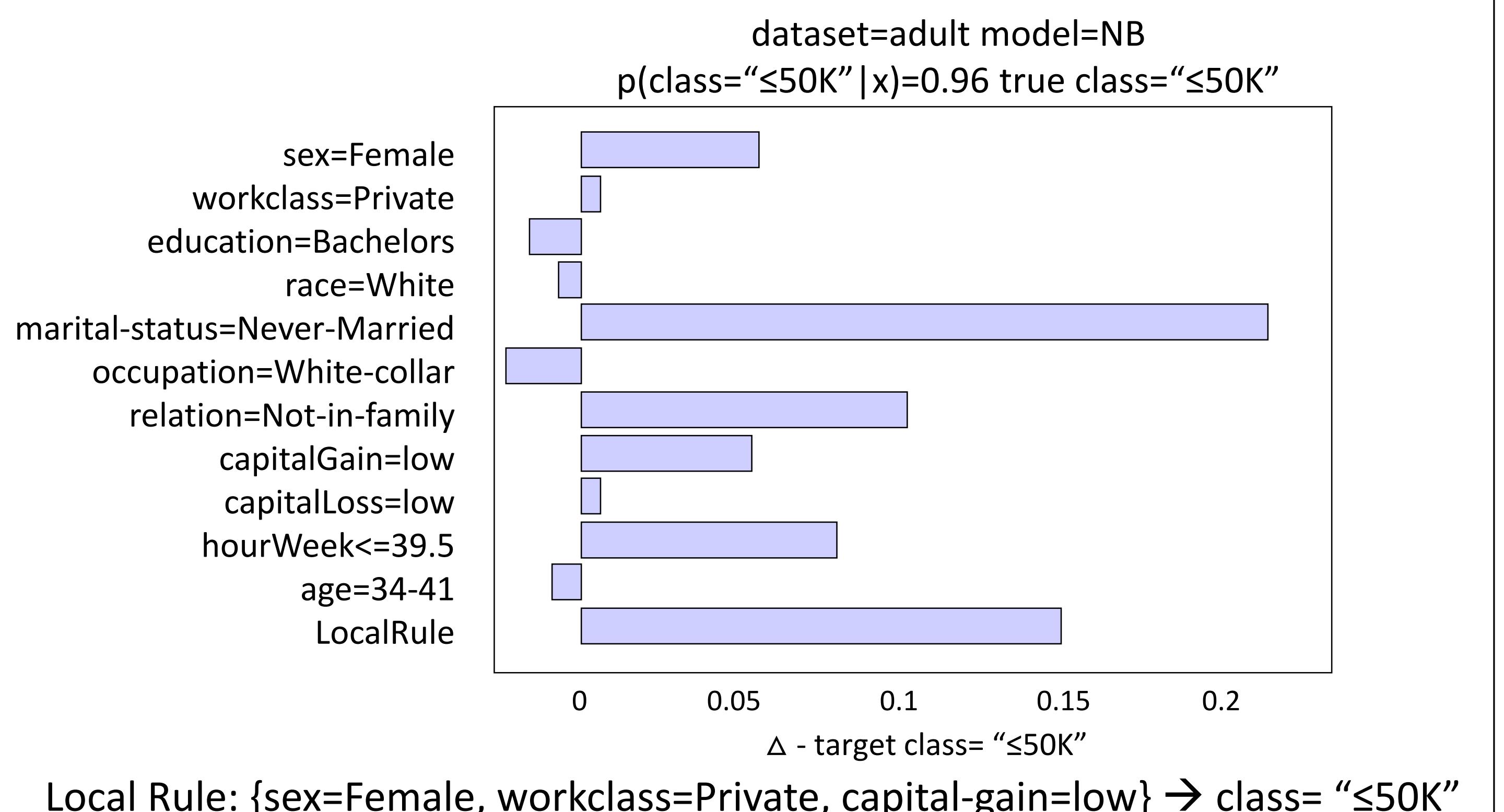
- Qualitative insight → **local rules**
Rules learned in the locality of the prediction
- Quantitative insight → **prediction difference**
Prediction change when one or more attribute values are omitted.

LACE Explanation Results

Explanation comparison



Insight of reasons behind predictions



Future work

- User studies → to assess LACE ability to provide actionable insights on model behavior
- Explanation of Big Data model predictions