

CALIME: Causality-Aware Local Interpretable Model-Agnostic Explanations



SCAN ME

Martina Cinquini, Riccardo Guidotti

Computer Science Department, University of Pisa, Italy

OVERVIEW

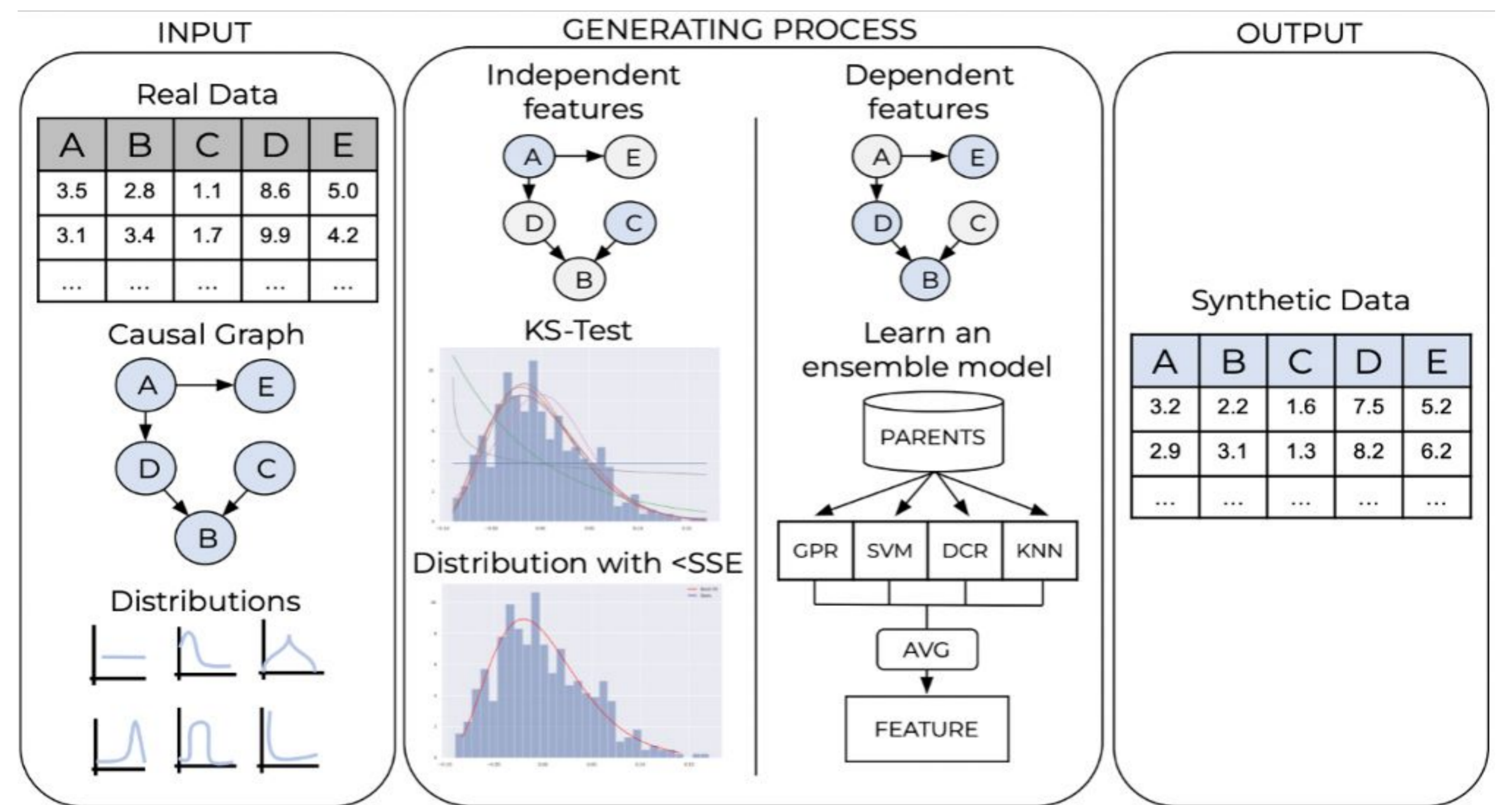
MOTIVATION

XAI approaches **do not take into account** causal relations among input features

OUR CONTRIBUTION

LIME variant **incorporates causal links** in the explanation extraction process

GENCDA



WHY DO WE NEED CAUSALITY?

Goal: Can the customer get the loan?

Black Box Prediction: No, the loan is denied

LIME Explanation: Low education level is mainly responsible

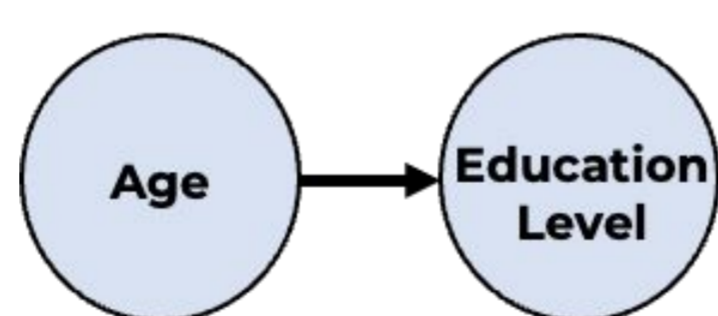
| Age | Income | Education Level | Weekly working hours |
|-----|--------|-----------------|----------------------|
| 24 | 800 | High School | 20 |
| 28 | 1300 | Bachelor Degree | 35 |
| ... | ... | ... | ... |

Generated Neighborhood

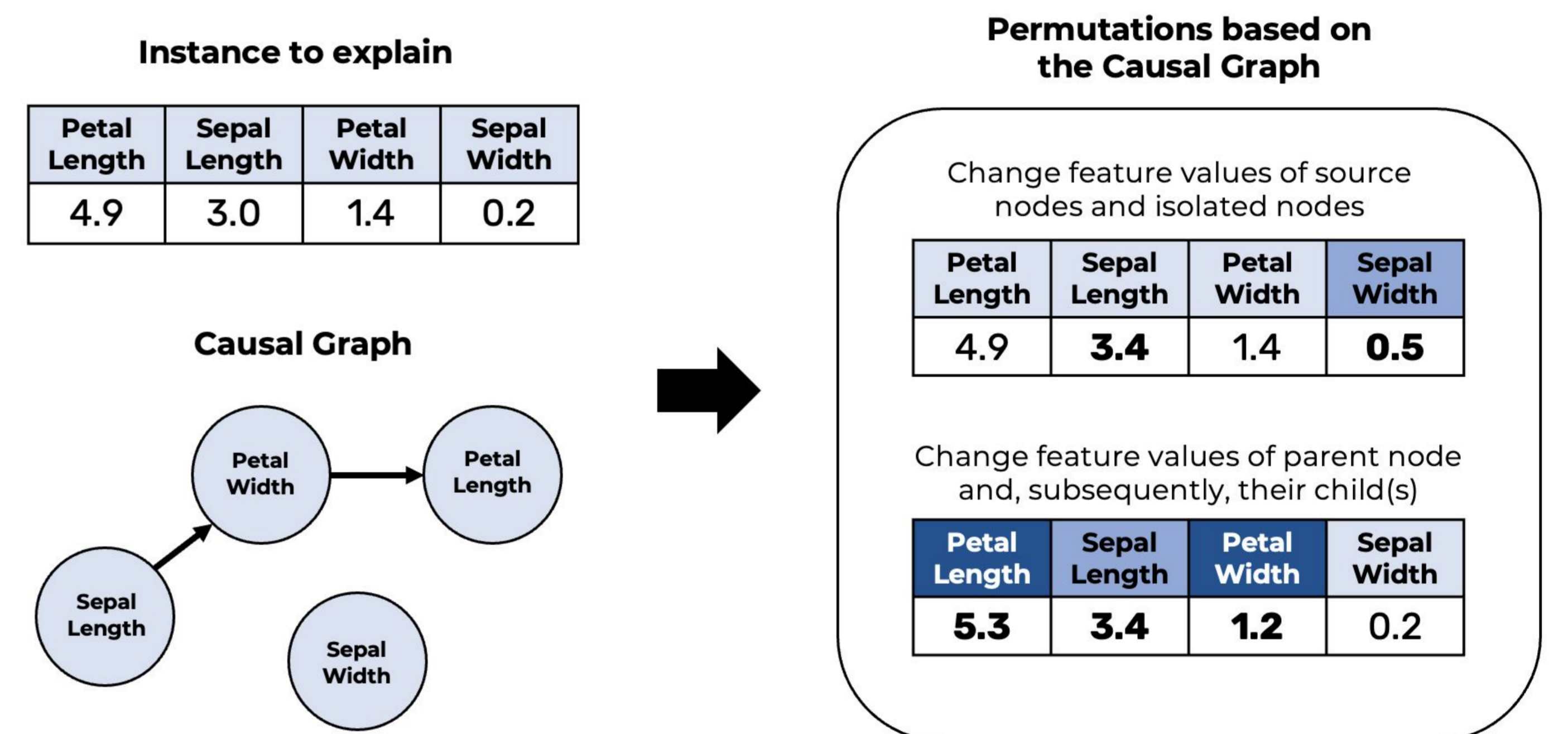
| | | | |
|----|-----|-----|----|
| 24 | 800 | PHD | 20 |
|----|-----|-----|----|



The generated instance is not plausible



NEIGHBORHOOD GENERATION



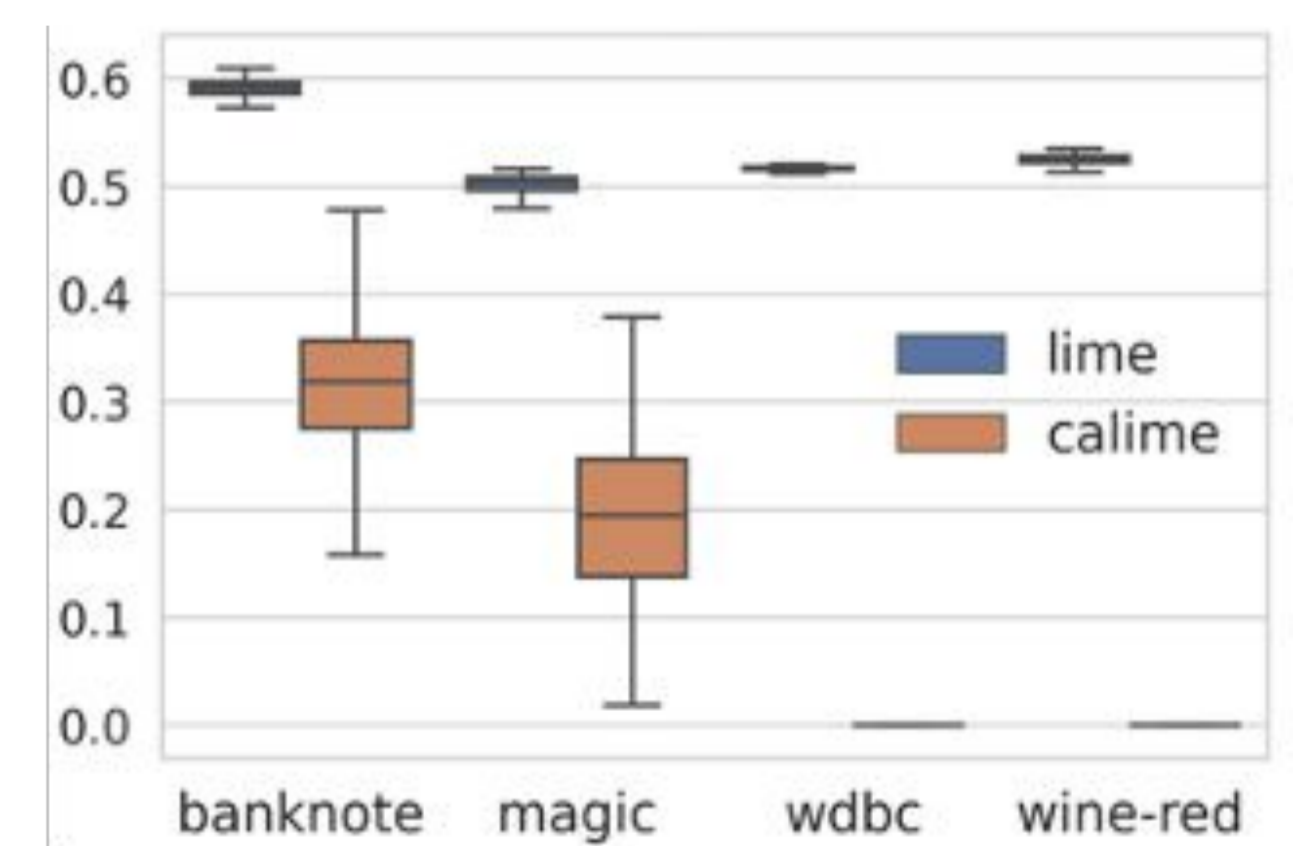
EXPERIMENTS

Datasets

banknote, wdbc, magic, wine-red

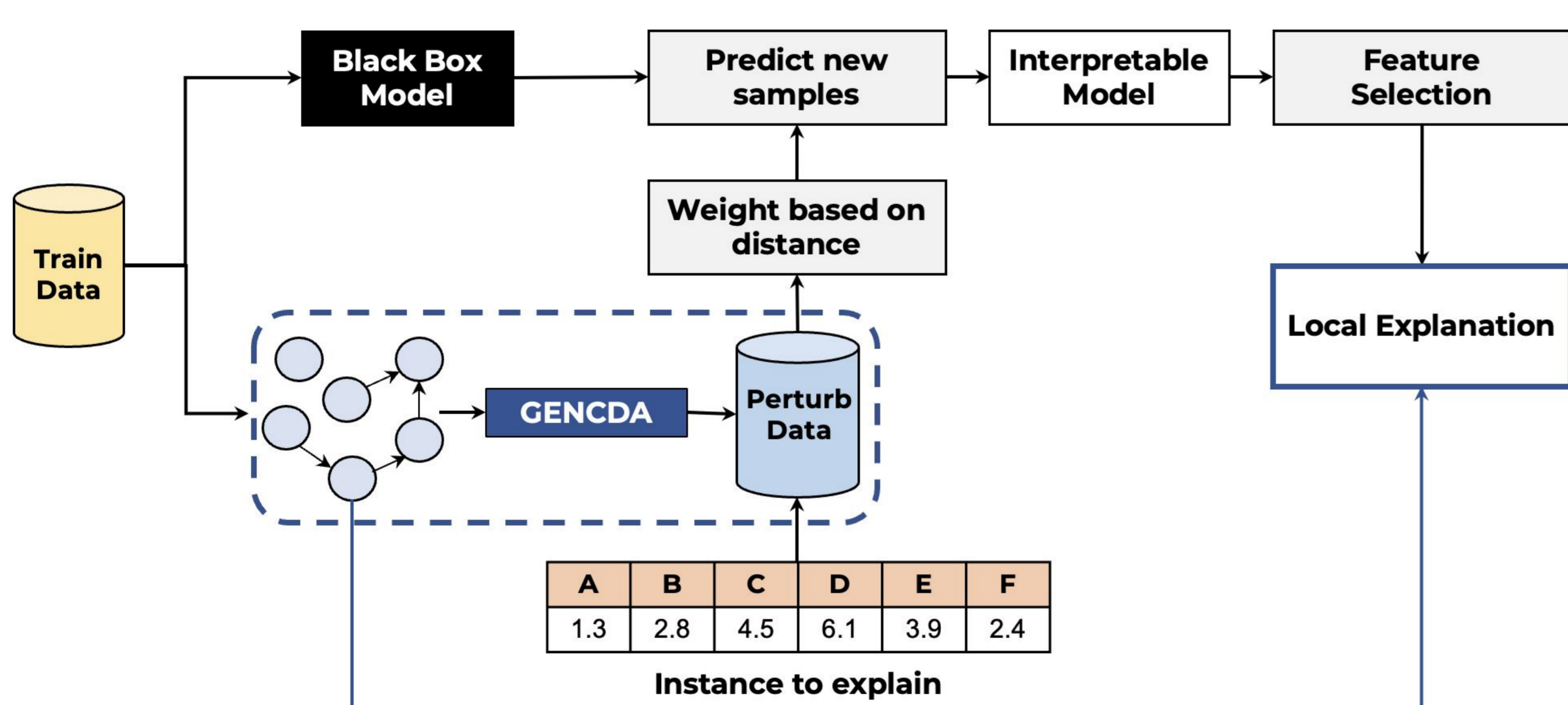
Evaluation metrics

Fidelity, Plausibility, Stability



CALIME outperforms LIME in both **black-box fidelity** and **explanations plausibility**

OUR FRAMEWORK



KEY TAKEAWAY

CALIME is the **first approach** able to **infer** and **integrate causal relations** to **promote interpretability** of Machine Learning models

