**Fairness** is not just about **data:** two models trained on the same data may be **fair** or **may not**.

We propose a **new way** to **analyze** and **measure** algorithmic fairness that goes **beyond predictive performance**.

### IFAFAM: Individual Feature Algorithmic Fairness Analysis Method applied on predicting students' performance

Mélina Verger<sup>1</sup>, Vanda Luengo<sup>1</sup>, Francois Bouchet<sup>1</sup>, Sébastien Lallé<sup>1</sup> <sup>1</sup> MOCAH team, Sorbonne University, CNRS, LIP6 image: Constant Structure Struct

### Motivation

Increasing need for algorithmic fairness analysis and quantification

Lack of **sensitive features discovery** – from what is learned from the models

Proliferation of *performance-oriented* fairness metrics and *performance-oriented* analysis approaches

# IFAFAM (contribution)

**General necessary steps** \* prediction probabilities (PP) retrieval \* sensitive-feature subset reduction

#### Qualitative analysis \* smoothing of PP distributions \* visual identification of unequal treatment

and stereotypical judgement

Quantitative analysis

\* a new way to measure algorithmic unfairness: the mean absolute density distance (MADD)

#### MADD (MAE-like but on density, model-independent)

$$rac{1}{N}\sum_{i=1}^{N}|d_i-d_i'|$$

*N*: total number of density observations  $d_i$ : density value related to the prediction probability  $p_i$  of one group  $d'_i$ : same as  $d_i$  but for the other group



#### Experimental validation :

- \* on several educational datasets and several binary classifier models (cross-tables) that predict students' success to courses
- \* comparison with existing fairness metrics (ABROCA, DI, TPR, ...)

# Discussion

- \* binary sensitive features
- \* supervised-learning oriented
  \* for any tabular data without preprocessing



1

2

3









