

On the Fairness ROAD: Robust Optimization for Adversarial Debiasing

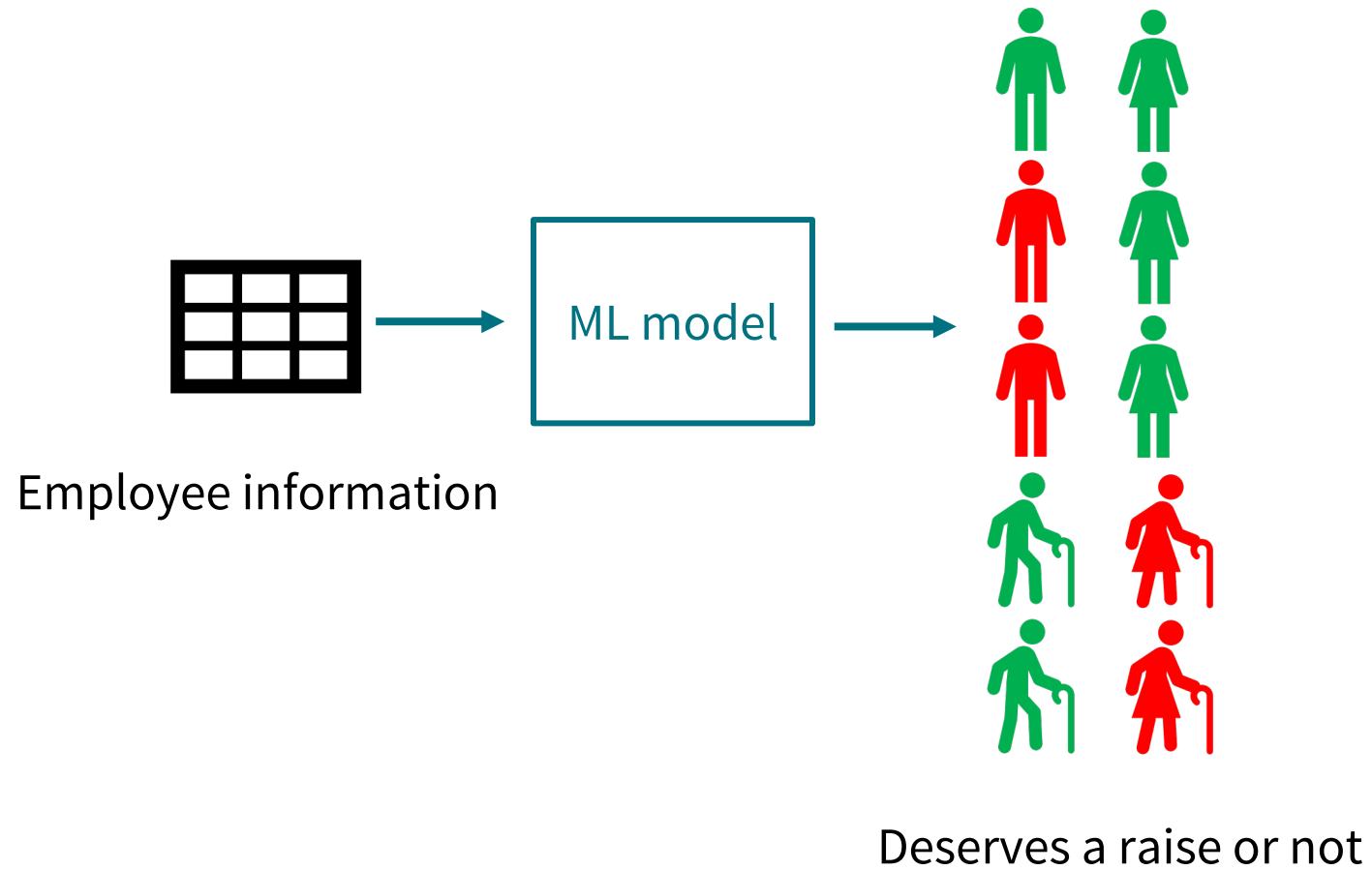
Vincent Grari*, Thibault Laugel*, Tatsunori Hashimoto, Sylvain Lamprier and Marcin Detyniecki



*Equal contribution



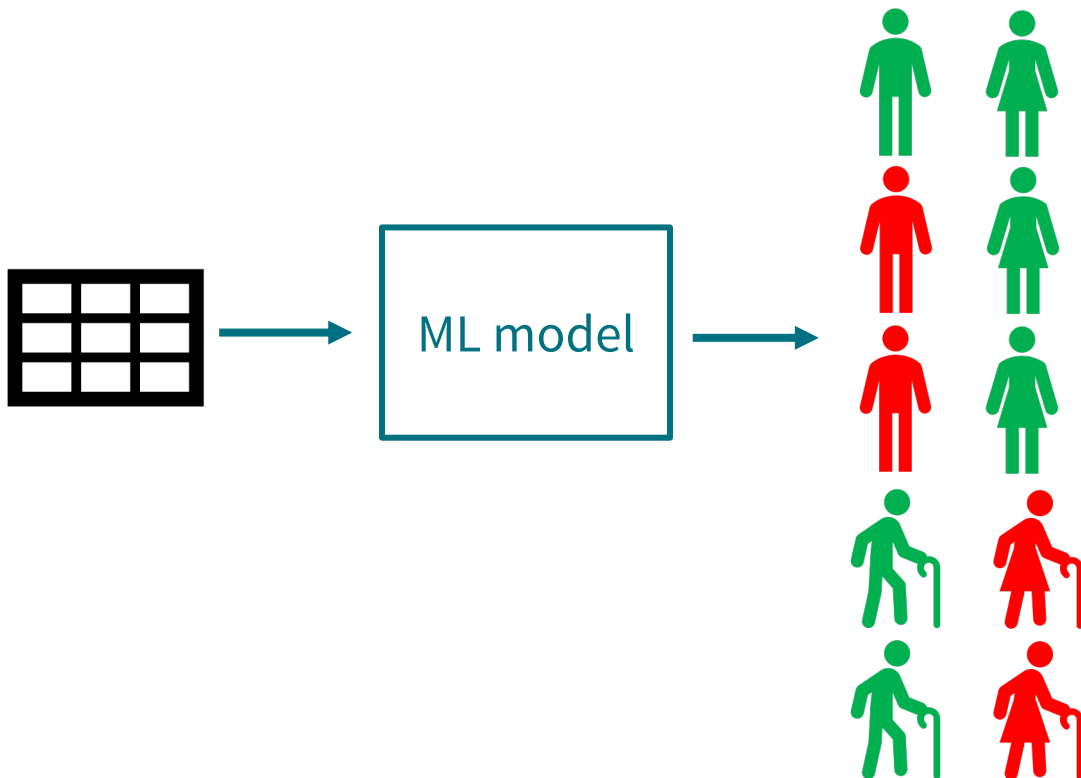
Context: algorithmic group fairness



Context: algorithmic group fairness

Traditional group fairness

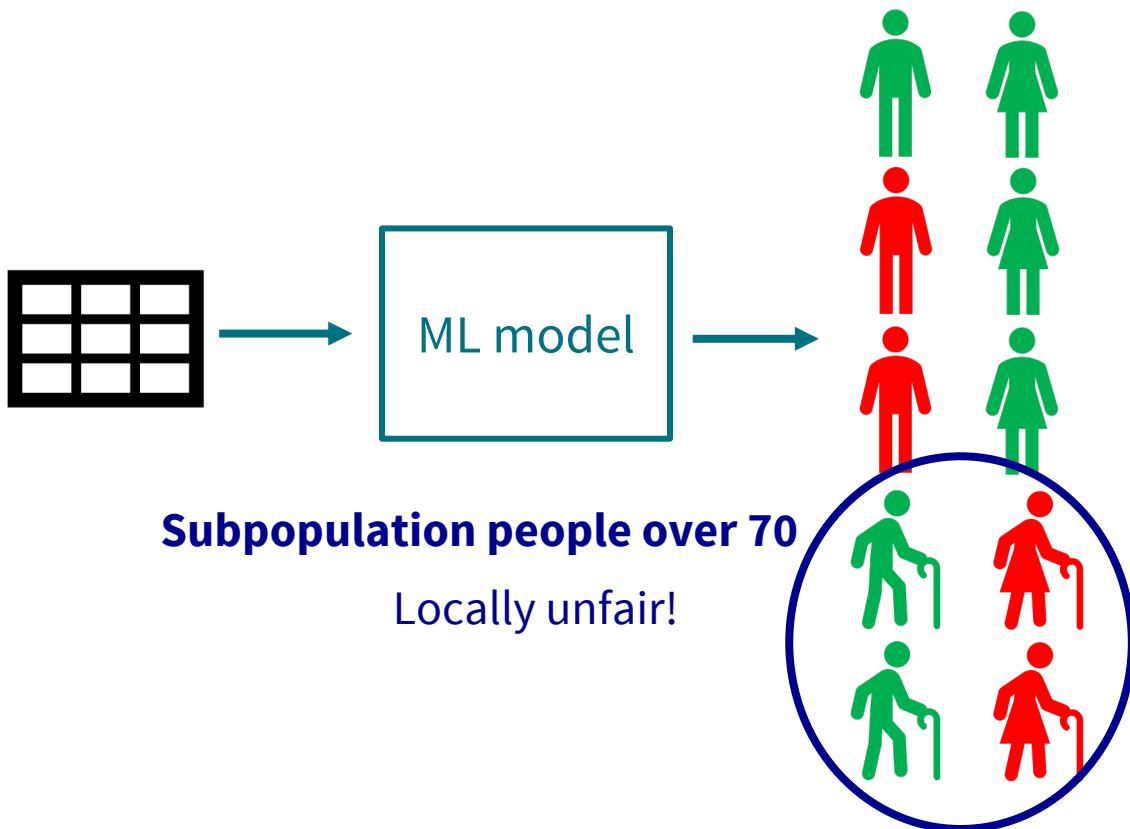
Globally fair model (DP): $\mathbb{P}(\hat{Y} = 1 | S = 1) = \mathbb{P}(\hat{Y} = 1 | S = 0)$



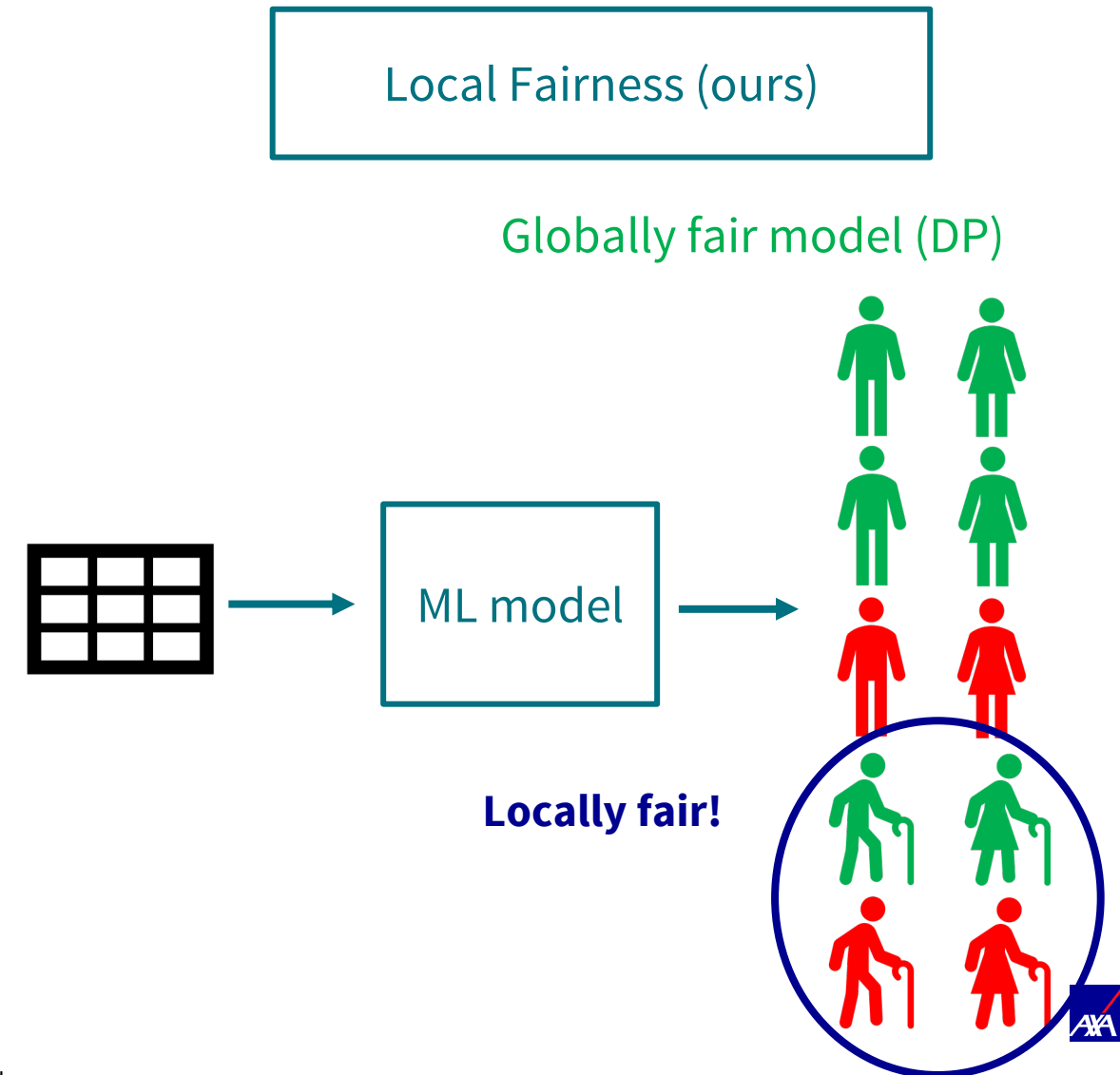
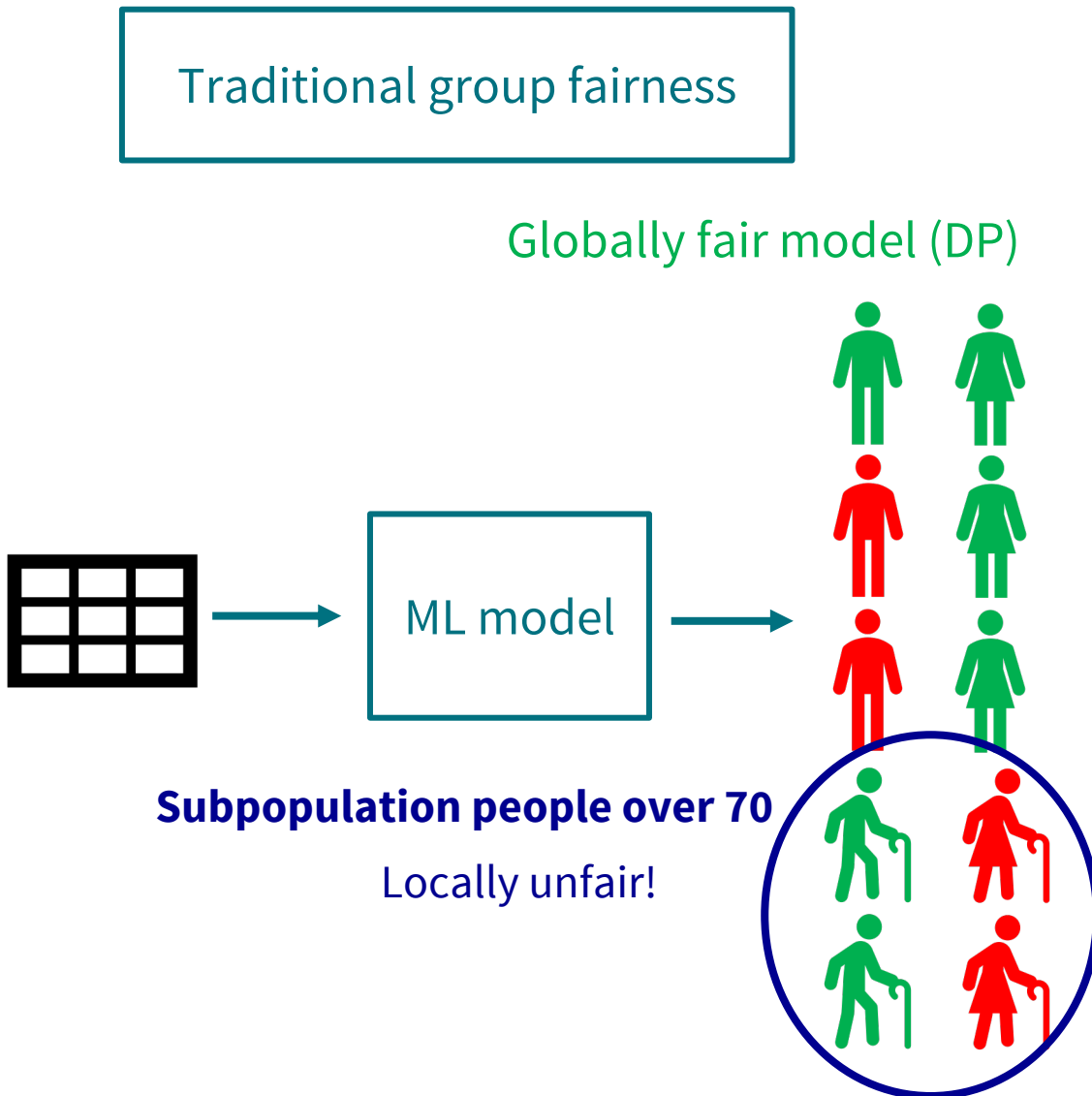
The Local (un)fairness problem

Traditional group fairness

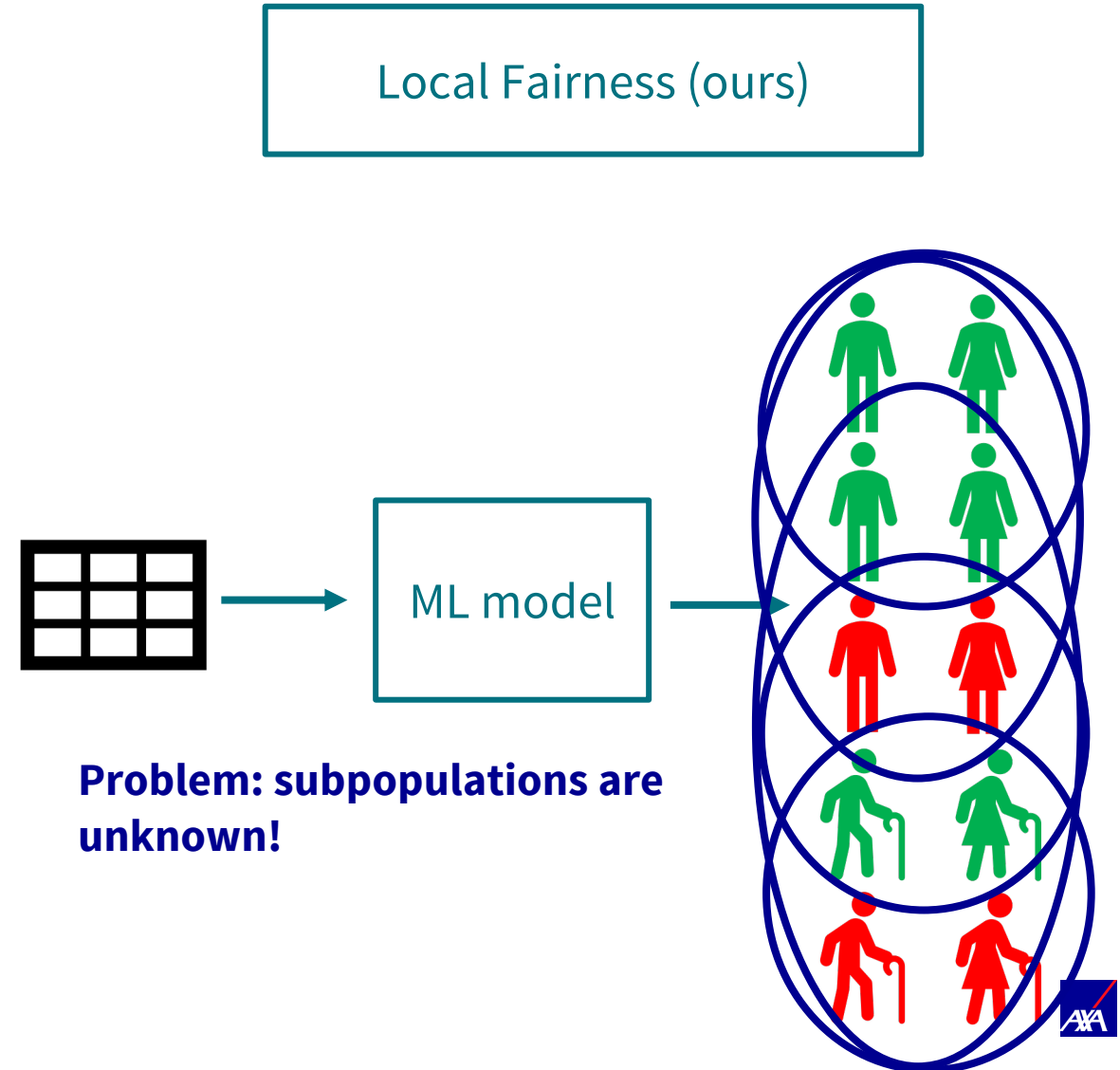
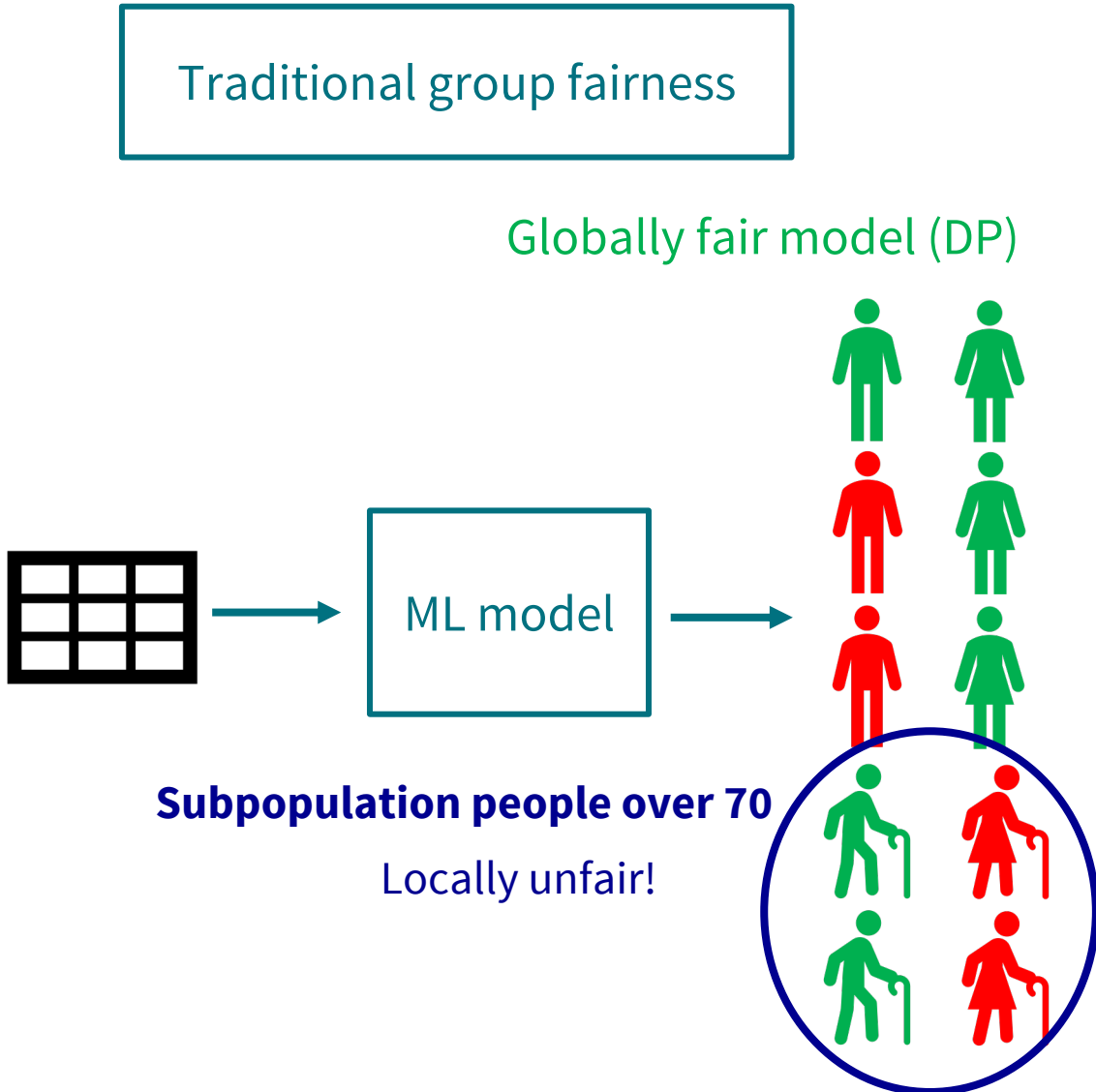
Globally fair model (DP)



The Local (Un)fairness problem



The Local (Un)fairness problem



Distributionally Robust Optimization (DRO) for Fairness

Traditional group fairness

$$\begin{aligned} \min_{w_f} \mathbb{E}_p[L_Y(f_{w_f}(x), y)] \\ \text{s. t. } DI_{(x,s) \sim p}(f_{w_f}(x), s) < \epsilon \end{aligned}$$

Local Fairness (ours)

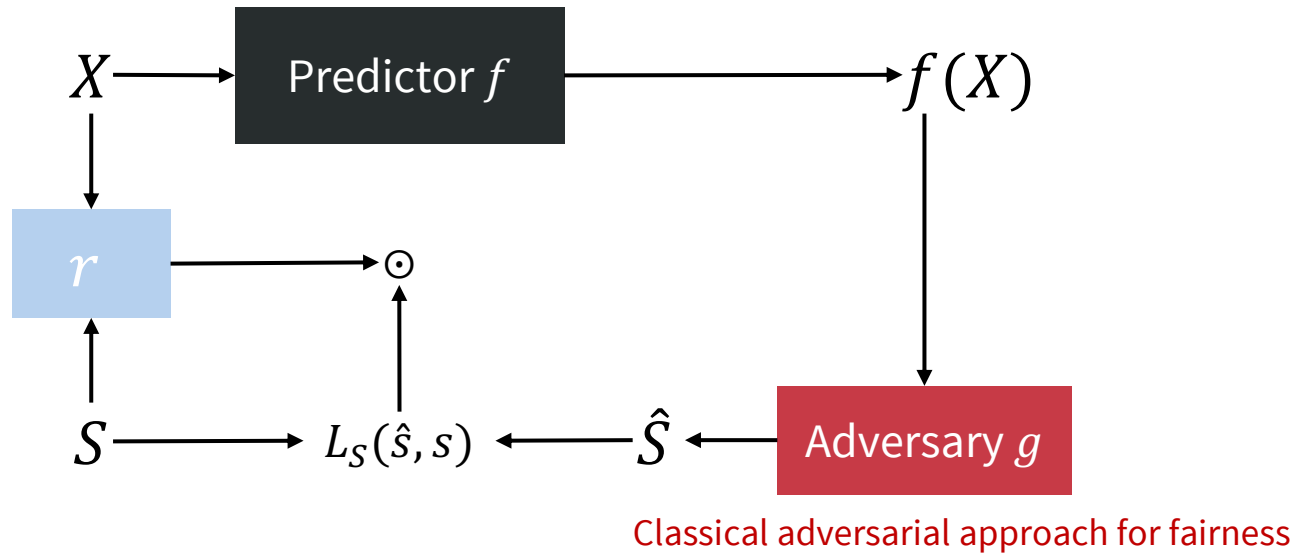
$$\begin{aligned} \min_{w_f} \mathbb{E}_p[L_Y(f_{w_f}(x), y)] \\ \text{s. t. } \max_{q \in Q} DI_{(x,s) \sim q}(f_{w_f}(x), s) < \epsilon \end{aligned}$$

Q : set of "plausible" distributions
~set of subpopulations

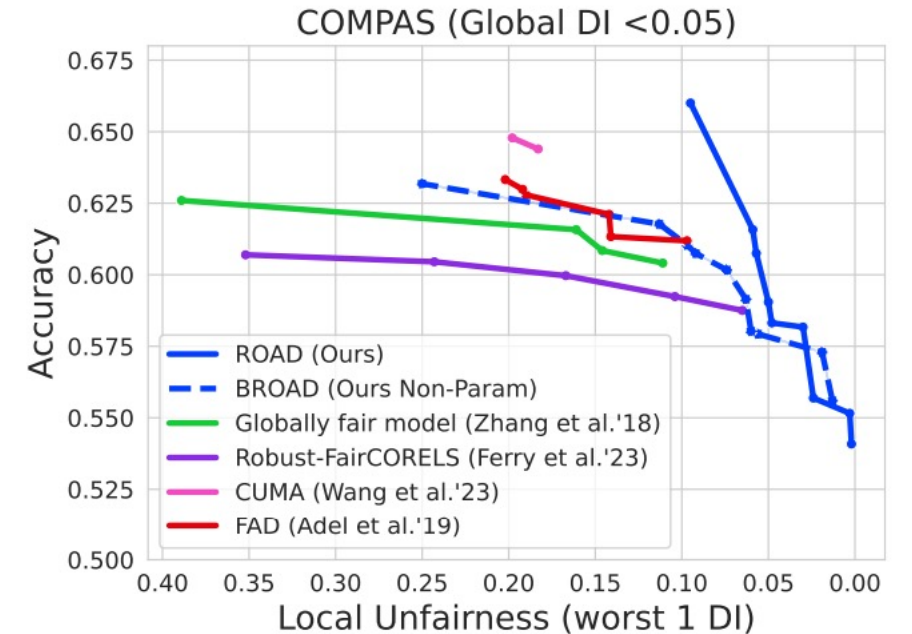
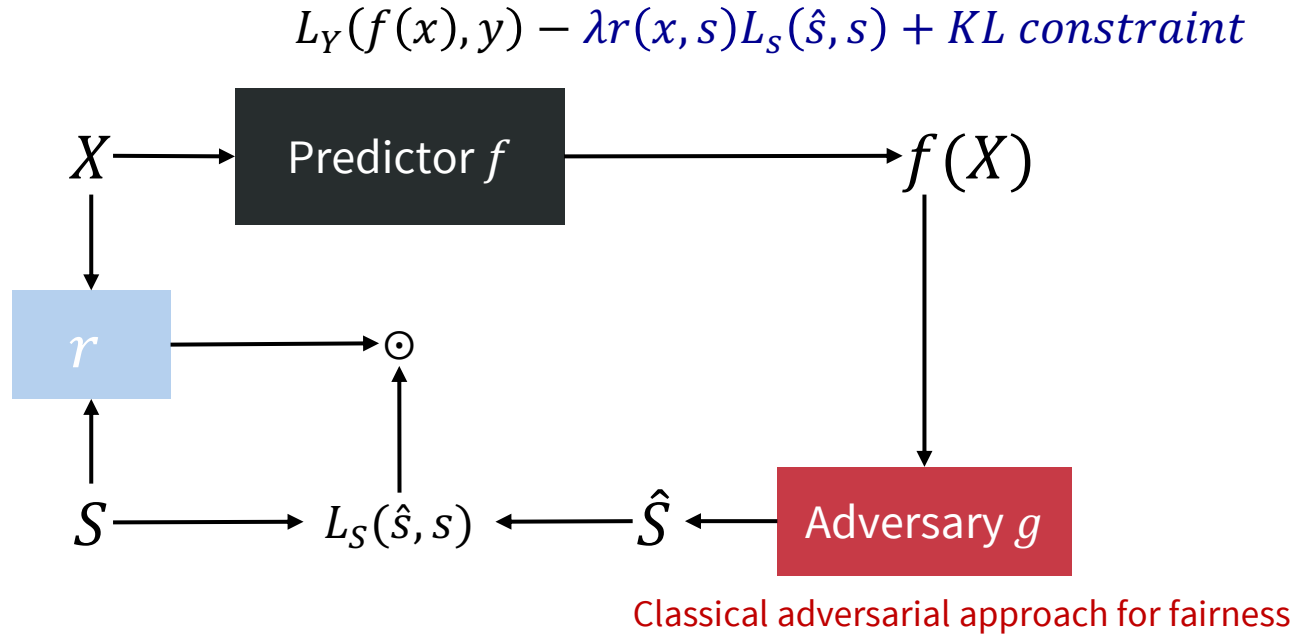
In practice: KL divergence-ball around p

Adversarial model for Distributionally Robust Fairness

$$L_Y(f(x), y) - \lambda r(x, s)L_S(\hat{s}, s) + KL \text{ constraint}$$



Adversarial model for Distributionally Robust Fairness



Results: more fair locally for the same levels of group fairness and accuracy

Thanks for watching!

Paper : <https://openreview.net/forum?id=xnhvVtZtLD>

Code: <https://github.com/axa-rev-research/ROAD-fairness/>