

Fair Wrapping for Black-box Predictions

Objectives

Exploit the properties (im)proper loss functions to develop a black-box classification post-processing method for algorithmic fairness.

We also want an approach with: \hookrightarrow flexibility to fairness criteria \hookleftarrow limited distortion of black-box \hookrightarrow algorithmic guarantees \hookleftarrow interpretability of correction

Black-box Post-Processing Setting

- ▶ \mathcal{X} a domain of observations;
 - ▶ $\mathcal{Y} \in \{-1, +1\}$ labels to classify;
 - ▶ S sensitive modalities partitioning \mathcal{X} such as age, race, etc.;
 - ▶ η_u an accurate but unfair black-box classifier.
- The goal of post-processing is to find a correction
- $$\eta_u \mapsto \eta_f \quad (1)$$
- which can be fitted to make η_f fairer than η_u .

Twist-Properness to Corrections

Class probability estimation loss functions define what the “best responses” of η is via

$$t_\ell(\eta) = \arg \inf_{\hat{\eta} \in [0,1]} L(\hat{\eta}, \eta). \quad (2)$$

The family of *improper* α -loss provides an (twist-proper) implementation of (1) through t_ℓ (2):

$$\eta_f(\mathbf{x}) = \frac{\eta_u(\mathbf{x})^{\alpha(\mathbf{x})}}{\eta_u(\mathbf{x})^{\alpha(\mathbf{x})} + (1 - \eta_u(\mathbf{x}))^{\alpha(\mathbf{x})}}. \quad (3)$$

TOPDOWN to Learn Alpha-Trees

With (3), our goal is to learn $\alpha(\cdot)$ to improve fairness.

To learn $\alpha(\cdot)$, we utilize a TOPDOWN tree induction algorithm for general risk minimization:

$$\eta_f \stackrel{\text{TOPDOWN}}{\longleftarrow} \min_{\alpha\text{-tree}} \mathbb{E}_{X \sim M_t} [L(\eta_f(X), \eta_t(X))] \quad (4)$$

- ▶ $\alpha(\cdot)$ is learnt as a *tree* (see Fig \checkmark);
- ▶ L is chosen to be the log-loss;
- ▶ M_t, η_t can be chosen for various fairness metrics.

Limitation on Distortion

Theorem 1. Suppose black-box η_u is bounded from extremes, s.t. $\text{Im}(\eta_u) \subseteq [(1 + \exp(B))^{-1}, (1 + \exp(-B))^{-1}]$, $0 < B \leq 3$ (Assum. 1)

and $|\alpha(\mathbf{x}) - 1| \leq 1/B$ (a.s.),

Then for any M : $\text{KL}(\eta_u, \eta_f; M) \leq \frac{\pi^2}{6 \cdot (2 + \exp(B) + \exp(-B))}$.

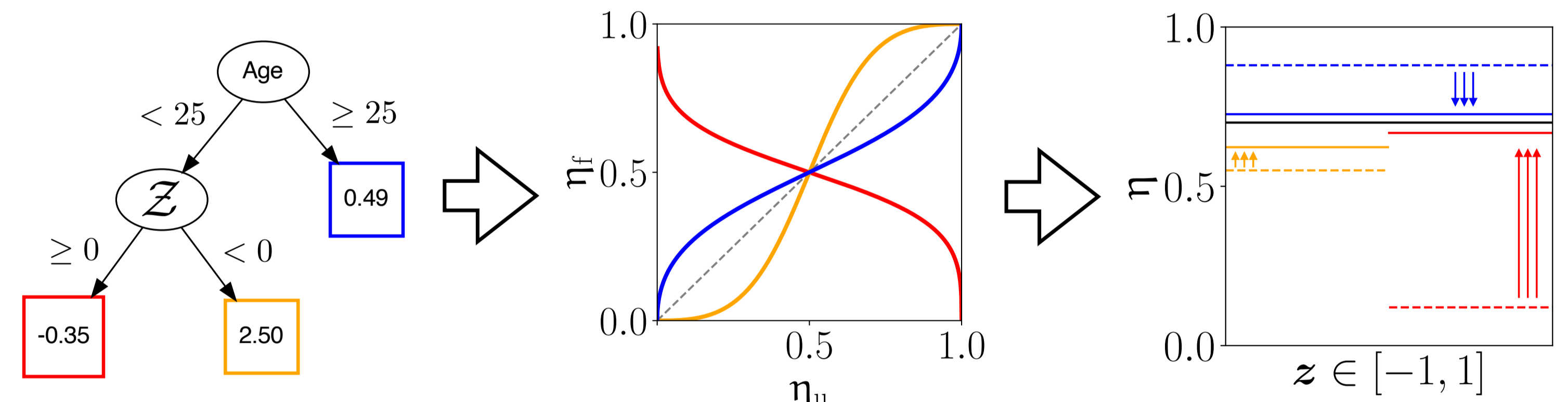


Figure: Toy Hiring Dataset: An example of an α -tree improving the fairness of a black-box (dotted colored) through a corrected posterior (solid colored).

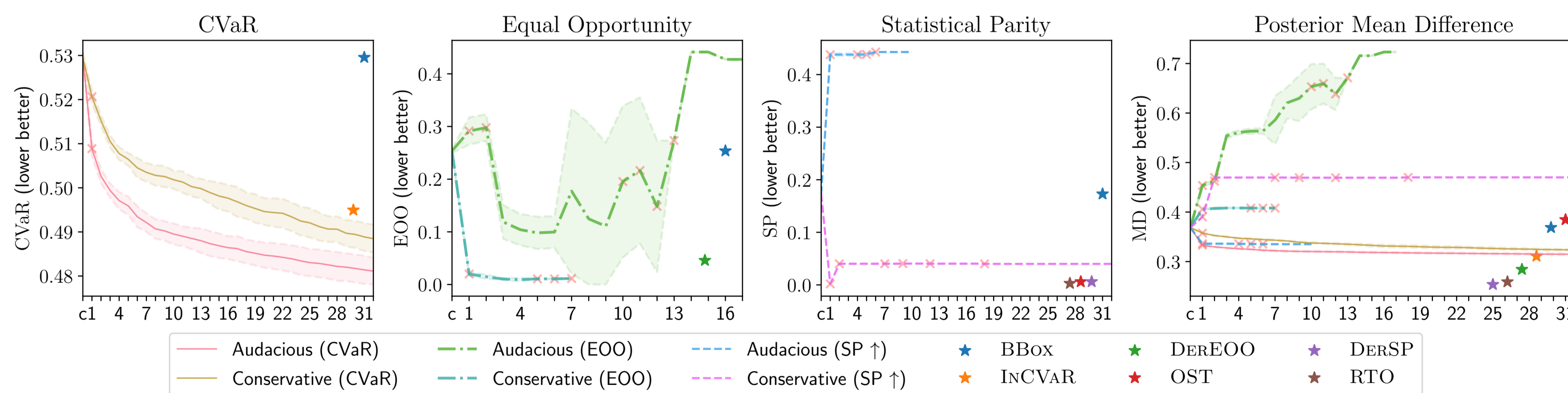


Figure: Improving fairness for a random forest black-box over a CA Income dataset for various criteria.

Algorithmic Convergence

Theorem 2. Suppose (Assum. 1) and a *Weak Learning Assumption* (with γ -witness) holds for all splits, then for any $\varepsilon > 0$, with TOPDOWN we have $\mathbb{E}_{X \sim M_t} [L(\eta_f(X), \eta_t(X))] \leq \varepsilon$ if

$$(\text{“\# Leaves in } \alpha\text{-tree”}) \quad |\Lambda(\alpha)| \geq (1/\varepsilon)^{c \log(1/\varepsilon)/\gamma^2}$$

- ▶ As this holds for any M_t, η_t , careful selection of targets allows for fairness guarantees for different criteria.