

Problem setup and fairness criterion

Observe R noisy crowd labels

$$Y_{1:R} = (Y_1, \dots, Y_R) \in \{0, 1\}^R$$

for a latent binary label $Y \in \{0, 1\}$, features $X \in \mathcal{X}$, and sensitive attribute $A \in \{0, 1\}$.

The **demographic-parity gap** of an aggregation rule ϕ is

$$\Delta_{\text{DP}}(Y^\phi) := |\mathbb{P}(Y^\phi = 1 \mid A = 1) - \mathbb{P}(Y^\phi = 1 \mid A = 0)|.$$

An aggregation rule ϕ is ε -fair when

$$\Delta_{\text{DP}}(Y^\phi) \leq \varepsilon.$$

Two reference aggregation rules

Majority vote

$$\mathbb{1}\left\{\sum_{r=1}^R Y_r \geq \frac{R}{2}\right\}.$$

Uses only the hard crowd labels.

Bayesian vote

$$\mathbb{1}\left\{\mathbb{P}(Y = 1 \mid Y_{1:R}, X, A) \geq \frac{1}{2}\right\}.$$

Uses the full posterior information.

Takeaway. We compare how standard crowd aggregation rules distort or preserve demographic parity.

Non-asymptotic fairness inequality

For $\phi \in \{\phi^*, \phi^{\text{MV}}\}$ and any R :

$$|\Delta_{\text{DP}}(Y^\phi) - \Delta_{\text{DP}}(Y)| \leq \sum_{a \in \{0,1\}} \mathbb{E}[e^{-RK_\phi(a,X)} \mid A = a].$$

The right-hand side decays exponentially in R : the fairness gap of Y^ϕ converges to the bias already present in the latent truth Y .

ε -Demographic Parity Gap of the Majority vote

$$\Delta_{\text{DP}}(Y^{\text{MV}}) \leq \varepsilon_R \sum_{r=1}^R \Delta_{\text{DP}}(Y_r),$$

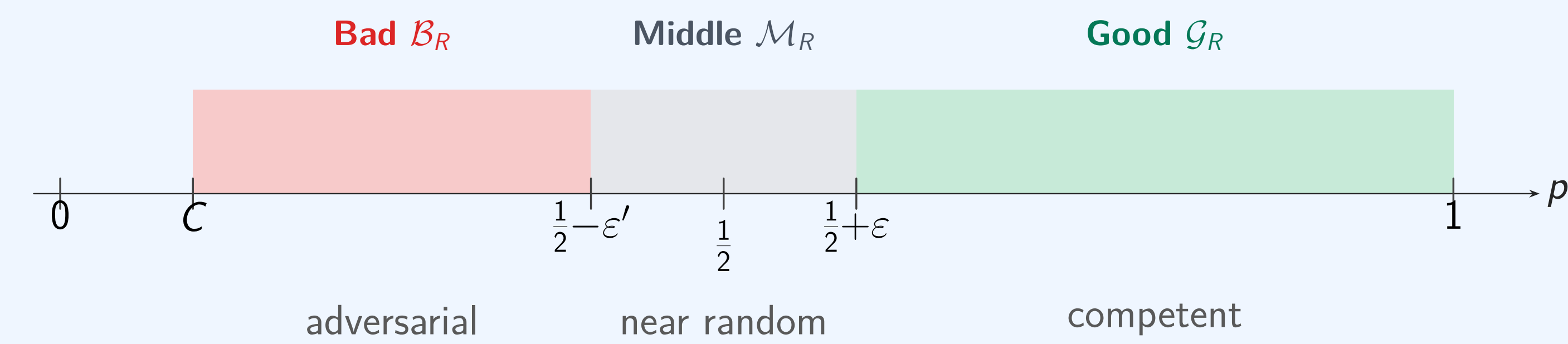
$$\varepsilon_R := \frac{\eta}{\min\{\sqrt{V_R(0)}, \sqrt{V_R(1)}\}},$$

with $\eta \simeq 0.4688$.

Interpretable crowdsourcing conditions

Rank annotators by accuracy

$$p_r := \mathbb{P}(Y_r = Y \mid X, A) \in [0, 1].$$



Majority vote:

$$\liminf_{R \rightarrow \infty} \left[\frac{|G_R|}{R} (1 + 2\varepsilon) + 2C \frac{|B_R|}{R} \right] > 1 \quad \text{a.s.}$$

Bayesian vote:

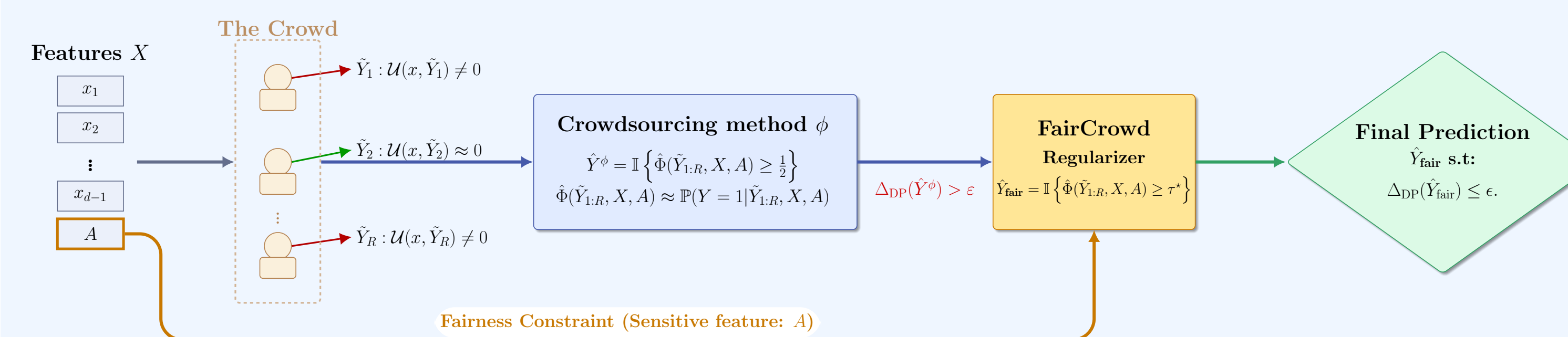
$$\sum_{r=1}^{\infty} \left(p_r - \frac{1}{2} \right)^2 = \infty \quad \text{a.s.}$$

Under these conditions:

$$\lim_{R \rightarrow \infty} \Delta_{\text{DP}}(Y_R^\phi) = \Delta_{\text{DP}}(Y), \quad \phi \in \{\phi^*, \phi^{\text{MV}}\}.$$

Takeaway. Large crowds recover, not erase, ground-truth bias.

FairCrowd: optimal ε -fair post-processing



Closed-form formula: boundary shift based on fairness trade-off

Instead of modifying the crowd model, FAIRCROWD solves

$$\min_{\phi \in \mathcal{G}_\varepsilon} \mathbb{P}(Y^\phi \neq Y) \quad \text{subject to} \quad \Delta_{\text{DP}}(Y^\phi) \leq \varepsilon.$$

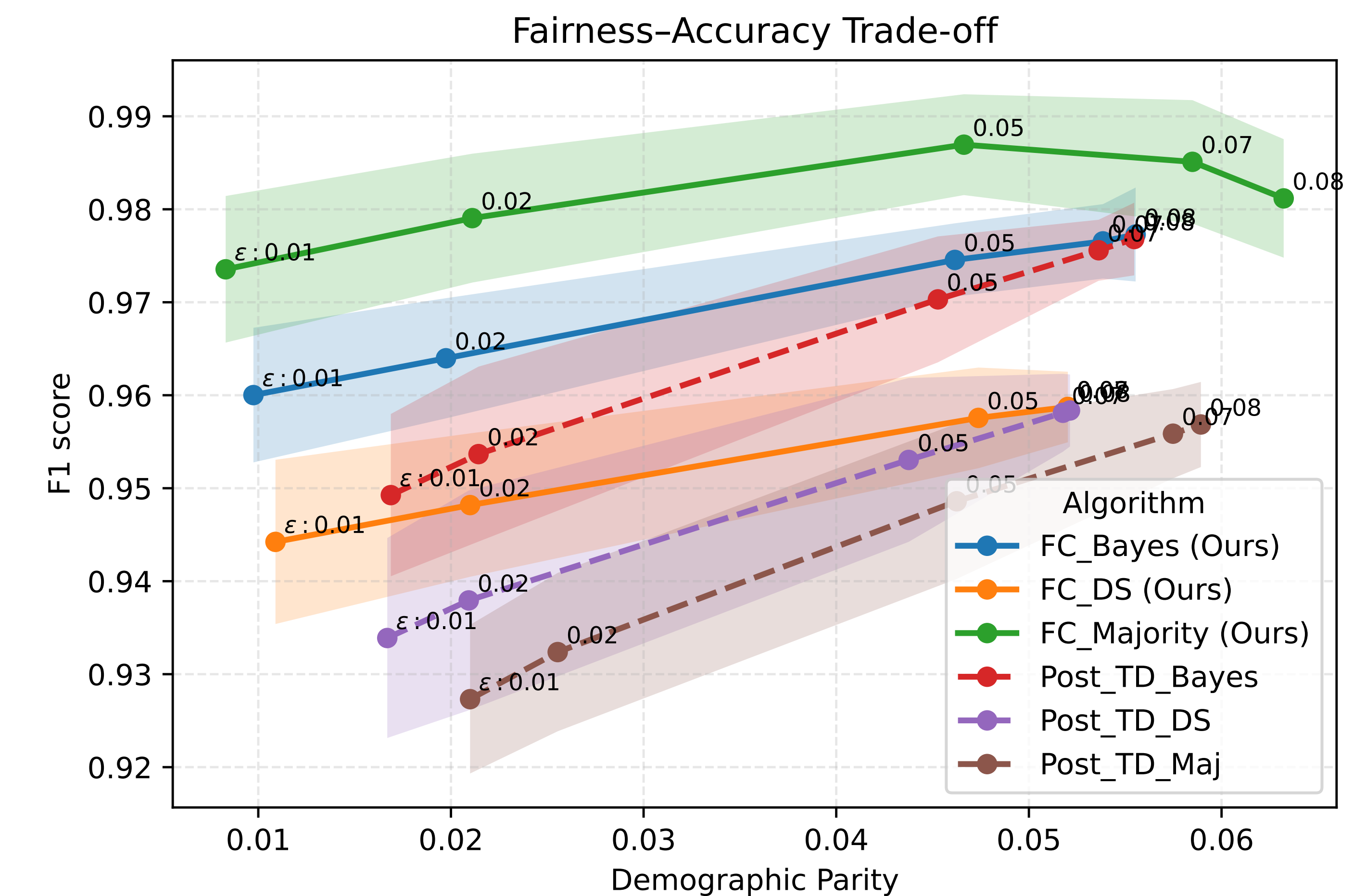
The closed-form optimal rule is a **group-dependent threshold shift**:

$$\phi_{\beta^*}^*(w, a) = \mathbb{1}\left\{\Phi_1(w, a) > \frac{\pi_a + (2a - 1)\beta^*}{2\pi_a}\right\},$$

where $\pi_a = \mathbb{P}(A = a)$ and β^* is the unique dual variable that saturates

$$\Delta_{\text{DP}}(Y^{\phi^*}) = \varepsilon.$$

Experiments: fairness–accuracy trade-off



Jigsaw Toxicity. Solid: FAIRCROWD. Dashed: Post TD baselines. Curves sweep ε .

Empirical results

Accuracy: FAIRCROWD enforces the prescribed fairness constraint while preserving, and in some cases improving, predictive performance compared to FairTD and Post TD baselines.

Scalability. On large datasets, the post-processing step remains computationally negligible, with a running time below 3 seconds.

Main references

- C. Denis, R. Elie, M. Hebiri, and F. Hu. *Fairness guarantees in multi-class classification with demographic parity*. JMLR, 25(130):1–46, 2024.
- C. Gao, Y. Lu, and D. Zhou. *Exact exponent in optimal rates for crowdsourcing*. ICML, pp. 603–611, 2016.
- S. Lazier, S. Thirumuruganathan, and H. Anahideh. *Fairness and bias in truth discovery algorithms*. arXiv:2304.12573, 2023.