

A Human-in-the-loop Framework to Construct Context-aware Mathematical Notions of Outcome Fairness

Mohammad Yaghini (University of Toronto), Andreas Krause (ETH Zürich) and Hoda Heidari (Carnegie Mellon University)

Motivation

Fairness is relative, complex and context-dependent

- No single off-the-shelf definition can capture it

How to define a context-aware notion of fairness?

- Bringing back human judgement into the decision-making loop

Two-pronged solution:

- Eliciting judgement by pairwise comparison
- Aggregating judgment through a social choice mechanism

Equality of Opportunity (EOP)

Distinguish morally-justifiable (**desert**) attributes from **circumstantial** ones.

A utility distribution CDF F under policy ϕ satisfies EOP if for all **circumstances** c, c' and all **desert levels** d

$$F^\phi(. | c, d) = F^\phi(. | c', d)$$

EOP Parameter Estimation

Circumstance. Estimate $\mathbf{c} := \mathbf{z}_p \in \mathbb{R}^k$ for participant p

To what extent do you agree with the following statement?
It is ethically acceptable for the attribute [...] to impact the decision a defendant receives.

Desert. Estimate $\mathbf{d}_p = \delta_p \cdot [x, y]$, where \mathbf{d}_p is not directly observable. Assuming there exists a linear function $D_p: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ such that $\mathbf{d}_p = D_p(x, y) = \delta_p \cdot [x, y]$, we wish to find δ_p

Desert Queries. Q pairwise questions about two scenarios $t_1^q = [x_1^q, y_1^q]$ and $t_2^q = [x_2^q, y_2^q]$

From an ethical standpoint, between the two decision subjects, who do you believe *deserves* a more lenient decision?

Utility. Similar to desert, estimate $u_p = v_p \cdot [x, y, \hat{y}]$ with Q pairwise questions about two scenarios $t_1^q = [x_1^q, y_1^q, \hat{y}_1^q]$ and $t_2^q = [x_2^q, y_2^q, \hat{y}_2^q]$:

..., who do you think *will benefit more* from their algorithmic decision?

MLE. Find δ_p and v_p that maximizes the likelihood of observed desert/utility differences

Preference Aggregation (Social Choice)

Borda Count.

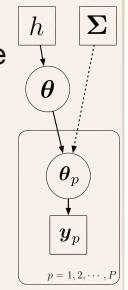
- Feature is circumstantial if most participants agree
- δ_p (and v_p) are averaged

Hierarchical Bayesian Model. Joint Parameter Estimation of θ : = society's preference vector

$\theta_p \sim \mathcal{N}(\theta, \Sigma)$ is participant p 's

$$\operatorname{argmin}_{\theta_p, \theta} - \sum_p \sum_q \log \Phi(a^q \theta_p \cdot [x_1^{p,q} - x_2^{p,q}, y_1^{p,q} - y_2^{p,q}])$$

s. t. $\|\theta_p - \theta\|_2 \leq \lambda, \|\theta_p\|_2 \leq 1, \|\theta\|_2 \leq 1$



Experimental Setup

From an ethical standpoint, between the two following decision subjects, who do you think will benefit more from their algorithmic decision?

ATTRIBUTE	SUBJECT #1	SUBJECT #2
Age Category	Older than 25	Older than 25
Race	Non-white	White
Gender	Male	Male
Charge Degree	Misdemeanor	Felony
Prior Counts	4	4
Algorithmic Decision	Low risk to reoffend	Low risk to reoffend
Actual Outcome	Will not reoffend	Will not reoffend

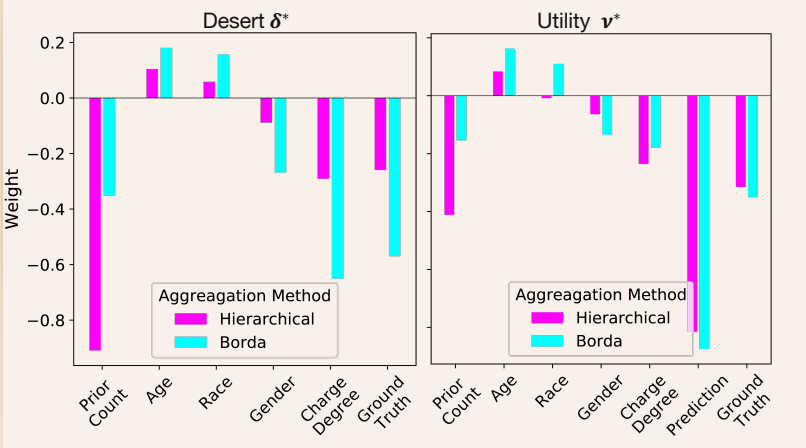
Proof-of-concept study ran on AMT

- 99 participants
- Conversational Interface

Note: The decision subject differences are marked in blue. If you are unsure about the meaning of any attribute, hold the cursor on it to see a definition.

Study Results (on AMT with 99 participants)

- Mixed opinion on age. Some thought younger people are less in control of their leniency



- EOP performs better in terms of equalizing the utility distributions (instead of pre-defined metrics) across groups
- EOP improves the utility for the whole populations; regardless of desert group

