



Group Retention when Using Machine Learning in Sequential Decision Making: the Interplay between User Dynamics and Fairness

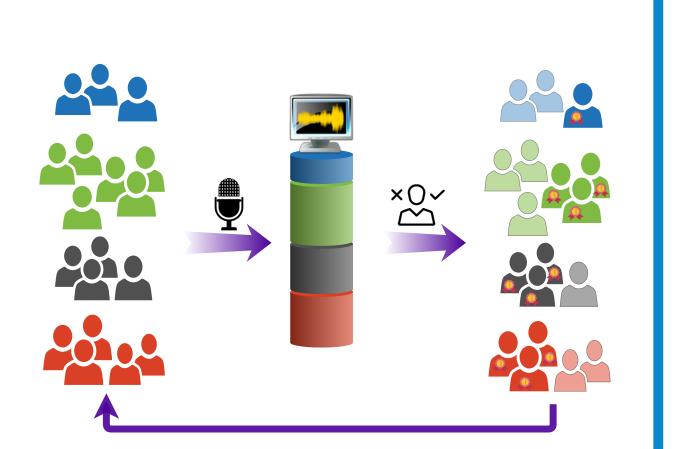
Paper

Xueru Zhang, Mohammad Mahdi Khalili, Cem Tekin and Mingyan Liu xueru@umich.edu; khalili@umich.edu; cemtekin@ee.bilkent.edu.tr; mingyan@umich.edu

OBJECTIVES

Motivation: ML model may be less favorable to groups contributing less to training process -> degrade population retention -> exacerbate representation disparity

- * What happens to the group representation over time when fair ML models are used.
- * How it is affected when underlying feature distributions are also reshaped by decisions.



PROBLEM FORMULATION

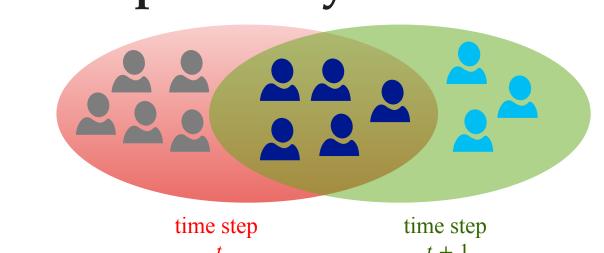
Two demographic groups G_a , G_b with time-varying X and $Y \in \{0, 1\}$

- feature distribution at t: $f_{k,t}(x) = g_{k,t}^0 f_{k,t}^0(x) + g_{k,t}^1 f_{k,t}^1(x)$, $x \in \mathbb{R}^d$
- representation disparity: $\frac{\overline{\alpha}_a(t)}{\overline{\alpha}_b(t)}$ with $\overline{\alpha}_a(t) + \overline{\alpha}_b(t) = 1$.

One-shot problem:

$$\min_{\theta_a,\theta_b} \ \mathcal{O}_t(\theta_a,\theta_b;\overline{\alpha}_a(t),\overline{\alpha}_b(t)) = \overline{\alpha}_a(t)O_{a,t}(\theta_a) + \overline{\alpha}_b(t)O_{b,t}(\theta_b)$$
s.t. $\Gamma_{\mathcal{C},t}(\theta_a,\theta_b) = 0$

Participation dynamics:



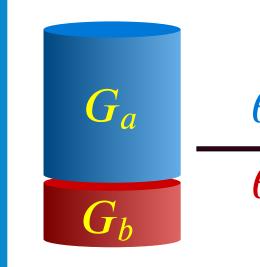
$$N_k(t+1) = N_k(t)\pi_{k,t}(\theta_k(t)) + \beta_k$$

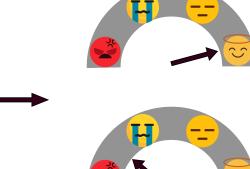
$$\overline{\alpha}_k(t+1) = \frac{N_k(t+1)}{N_k(t+1)}$$

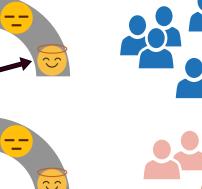
Goal: what happens to $\frac{\overline{\alpha}_a(t)}{\overline{\alpha}_b(t)}$ when one-shot fair decisions are applied $\forall t$.

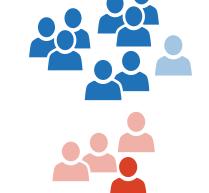
MONOTONICITY CONDITION (MC)

Two one-shot problems $\widehat{O}(\theta_a, \theta_b; \widehat{\alpha}_a, \widehat{\alpha}_b)$ and $\widetilde{O}(\theta_a, \theta_b; \widetilde{\alpha}_a, \widetilde{\alpha}_b)$ defined over $\widehat{f}_k(x)$ and $\widetilde{f}_k(x)$ satisfy MC under a dynamic model if $\forall \ \widetilde{\alpha}_a + \widetilde{\alpha}_b = 1$ and $\widehat{\alpha}_a + \widehat{\alpha}_b = 1$ such that $\frac{\widehat{\alpha}_a}{\widehat{\alpha}_b} < \frac{\widetilde{\alpha}_a}{\widetilde{\alpha}_b}$, the resulting retention rates satisfy $\widehat{\pi}_a(\widehat{\theta}_a) < 1$ $\widetilde{\pi}_a(\widetilde{\theta}_a)$ and $\widehat{\pi}_b(\widehat{\theta}_b) > \widetilde{\pi}_b(\widetilde{\theta}_b)$.

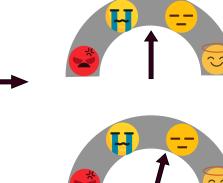




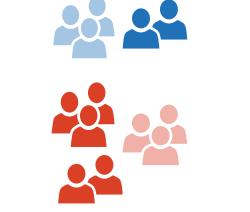


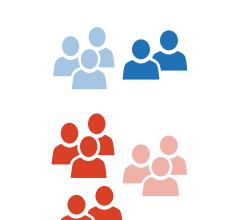








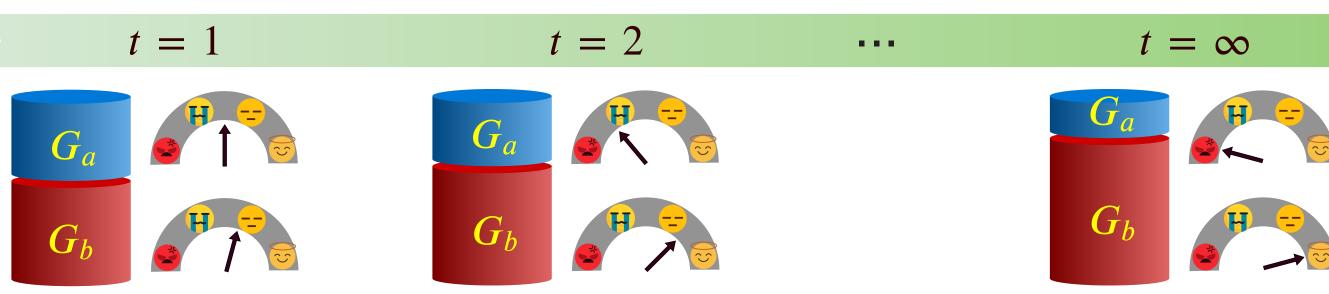




EXACERBATION OF REPRESENTATION DISPARITY

• Monotonic change of $\frac{\overline{\alpha}_a(t)}{\overline{\alpha}_b(t)}$ under MC:

Consider $\{O_t(\theta_a, \theta_b; \overline{\alpha}_a(t), \overline{\alpha}_b(t))\}_{t=1}^{\infty}$, if O_t and O_{t+1} satisfy MC $\forall t$ under a dynamic model, then we have: if $\pi_{a,1}(\theta_a(1)) < \pi_{b,1}(\theta_b(1))$ (different retention), then $\frac{\overline{\alpha}_a(t+1)}{\overline{\alpha}_b(t+1)} < \frac{\overline{\alpha}_a(t)}{\overline{\alpha}_b(t)}$ (exacerbation of disparity) and $\pi_{a,t+1}(\theta_a(t+1))$ 1)) $<\pi_{a,t}(\theta_a(t))$ $<\pi_{b,t}(\theta_b(t))$ $<\pi_{b,t+1}(\theta_b(t+1))$ (discrepancy increases).



CASES SATISFYING MC UNDER THE SAME $f_k(x)$

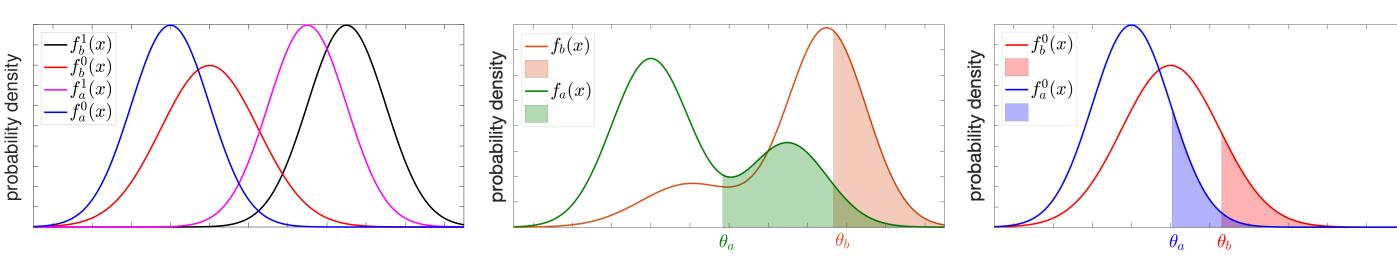
Two one-shot problems $\widehat{O}(\theta_a, \theta_b; \widehat{\alpha}_a, \widehat{\alpha}_b) = \widehat{\alpha}_a O_a(\theta_a) + \widehat{\alpha}_b O_b(\theta_b)$ and $\widetilde{O}(\theta_a, \theta_b; \widetilde{\alpha}_a, \widetilde{\alpha}_b) = \widetilde{\alpha}_a O_a(\theta_a) + \widetilde{\alpha}_b O_b(\theta_b)$ satisfy MC when:

- $*O_k(\widehat{\theta}_k) \neq O_k(\widehat{\theta}_k)$ for any possible $\widehat{\alpha}_k \neq \widetilde{\alpha}_k$
- * they are defined over the same $f_k(x)$
- * dynamics satisfy $\pi_k(\theta_k) = h_k(O_k(\theta_k))$ for some decreasing function $h_k(\cdot)$ ⇒ hold under commonly used objectives, dynamics and fairness criteria!

Technical Assumption:

A ONE-DIMENSIONAL THRESHOLD DECISION

- Decision rule:
 - $h_{\theta}(x) = \mathbf{1}(x \geq \theta), X \in \mathbb{R}$
- Objective function: $O_{k,t}(\theta_k) = L_{k,t}(\theta_k)$ with 0-1 loss $\mathbf{1}(y \neq h_{\theta}(x))$
- Fairness criterion $\Gamma_{\mathcal{C},t}(\theta_a,\theta_b)\longleftrightarrow \theta_a=\phi_{\mathcal{C},t}(\theta_b)$:
- Simple fair (Simple): $\Gamma_{\mathcal{C},t} = \theta_a \theta_b$
- Statistical Parity (StatPar): $\Gamma_{\mathcal{C},t} = \int_{\theta_a}^{\infty} f_{a,t}(x) dx \int_{\theta_b}^{\infty} f_{b,t}(x) dx$
- Equal Opportunity (EqOpt): $\Gamma_{\mathcal{C},t} = \int_{\theta_a}^{\infty} f_{a,t}^0(x) dx \int_{\theta_b}^{\infty} f_{b,t}^0(x) dx$



→ One-shot solutions under either Simple, EqOpt or StatPar:

- Bounded by a fixed interval $\forall \overline{\alpha}_a(t), \overline{\alpha}_b(t)$:
- $\theta_a(t) \in [\phi_{\mathcal{C},t}(\delta_{b,t}), \delta_{a,t}], \theta_b(t) \in [\delta_{b,t}, \phi_{\mathcal{C},t}^{-1}(\delta_{a,t})]; \quad \delta_{k,t} = \operatorname{argmin}_{\theta} L_{k,t}(\theta)$
- \exists a function $\Psi_{\mathcal{C},t}(\theta_a,\theta_b)$ increasing in θ_a,θ_b such that:

$$\Psi_{\mathcal{C},t}(\theta_a(t),\theta_b(t)) = \frac{\overline{\alpha}_a(t)}{\overline{\alpha}_b(t)}$$

RESULTS UNDER THE FIXED $f_k(x)$

- Dynamics:
- 1. By model accuracy: $\pi_k(\cdot) =
 u(L_k(\cdot)) ext{ with }
 u(\cdot) \!\!\downarrow \quad {}^{\scriptscriptstyle{0.00^{\perp}}}_{\scriptscriptstyle{0.04^{\perp}}}$
- 2. By intra-group disparity: $\pi_k(\cdot) = w(D_k(\cdot)) ext{ with } w(\cdot) \uparrow \circ_{\scriptscriptstyle{\mathsf{0.00}}}$
- **Exacerbation of representation disparity:**

Consider $\left\{ O_t(\theta_a, \theta_b; \overline{\alpha}_a(t), \overline{\alpha}_b(t)) \right\}_{t=1}^{\infty}$ under either Simple, EqOpt or StatPar criterion, when $f_{k,t}(x) = f_k(x), \forall t$:

- O_t and O_{t+1} satisfy MC $\forall t$ under above dynamics.
- $\{(\theta_a(t), \theta_b(t))\}_{t=1}^{\infty}$ converges monotonically to a constant decision.

RESULTS UNDER TIME-VARYING $f_{k,t}(x)$

- Dynamics driven by model accuracy: $\pi_k(\cdot) = \nu(L_k(\cdot))$
- Time-varying feature distributions: $f_{k,t}(x) = g_{k,t}^0 f_{k,t}^0(x) + g_{k,t}^1 f_{k,t}^1(x)$ G_k^0 and G_k^1 may react differently to the same θ_k :

 $\qquad \qquad \mathbf{f}_{k}^{j}(\mathbf{x}) \qquad \qquad \mathbf{g}_{k,t}^{j}\mathbf{f}_{k,t}^{j}(\mathbf{x}) \qquad \qquad \mathbf{g}_{k,t+1}^{j}\mathbf{f}_{k,t+1}^{j}(\mathbf{x})$

- (i) $f_{k,t}^j(x) = f_k^j(x), \forall t \text{ but } g_{k,t}^j$ changes: G_k^j 's retention is determined by $L_{k,t}^j$.
- (ii) $g_{k,t}^j = g_k^j, \forall t \text{ but } f_{k,t}^j(x) \text{ changes:}$
- For G_k^i that is less favored by decisions, its members may make extra effort for lowering their losses.
- => Exacerbation in representation disparity can accelerate:
- **(P1):** $\left\{ \boldsymbol{O}^{f}_{t}(\theta_{a}, \theta_{b}; \overline{\alpha}_{a}^{f}(t), \overline{\alpha}_{b}^{f}(t)) \right\}_{t=1}^{\infty} \text{ over } f_{k}(x) \text{ with } \pi_{a}^{f}(\theta_{a}^{f}(1)) < \pi_{b}^{f}(\theta_{b}^{f}(1))$ (P2): $\left\{ \boldsymbol{O}^{v}_{t}(\theta_{a}, \theta_{b}; \overline{\alpha}_{a}^{v}(t), \overline{\alpha}_{b}^{v}(t)) \right\}_{t=1}^{\infty} \text{ over } \left\{ f_{k,t}(x) \right\}_{t=1}^{\infty} \text{ with } f_{k,1}(x) = f_{k}(x)$ Under either Simple, EqOpt or StatPar, if $f_{k,t}(x)$ changes slowly w.r.t. the change in one-shot decisions:
 - O_t^v and O_{t+1}^v satisfy MC $\forall t$ under this dynamics.
 - $-\frac{\overline{\alpha}_a^v(t)}{\overline{\alpha}_b^v(t)} < \frac{\overline{\alpha}_a^f(t)}{\overline{\alpha}_b^f(t)}, \forall t$

MITIGATION WHEN $\pi_{k,t}(\cdot) = \nu(L_{k,t}(\cdot))$

Equalized Loss (EqLos) fairness:

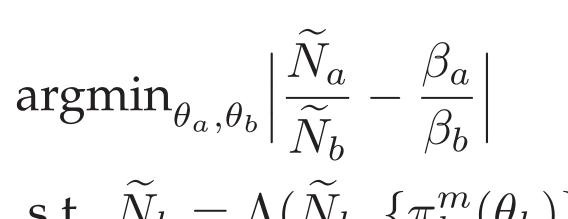
- $\Gamma_{\mathcal{C},t} = L_{a,t}(\theta_a) L_{b,t}(\theta_b)$ * It maintains group representation:
- $\lim_{t\to\infty} \frac{\overline{\alpha}_a(t)}{\overline{\alpha}_b(t)} = \frac{\beta_a}{\beta_b}$
- => Fairness has to be defined with a good understanding of underlying participation dynsmics!

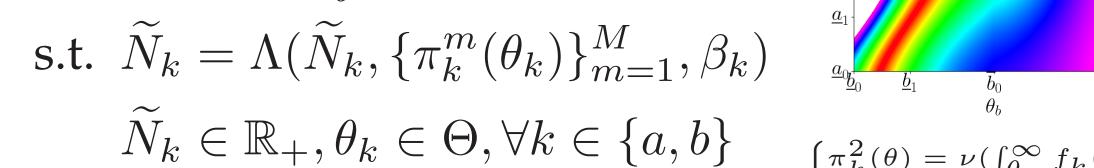
MITIGATION FOR GENERAL DYNAMICS

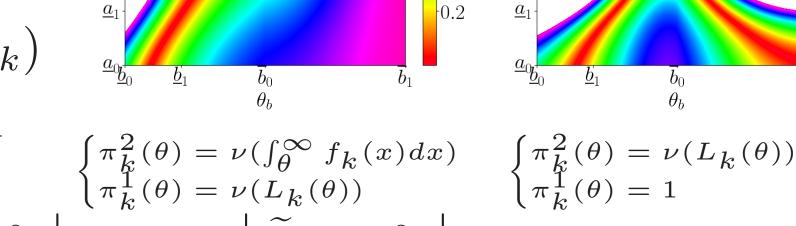
Find the proper criterion for some general dynamics when $f_{k,t}(x) =$ $f_k(x), \forall t$:

$$N_k(t+1) = \Lambda(N_k(t), \{\pi_k^m(\theta_k(t))\}_{m=1}^M, \beta_k)$$

• Assume $\exists (\theta_a, \theta_b)$ s.t. corresponding dynamics have stable fixed points.

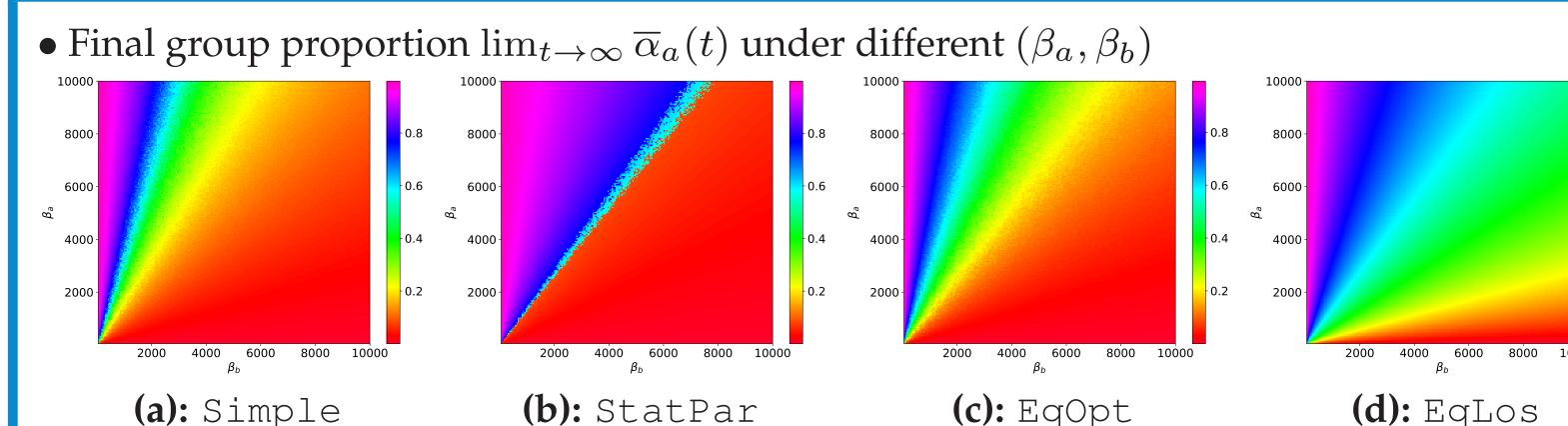


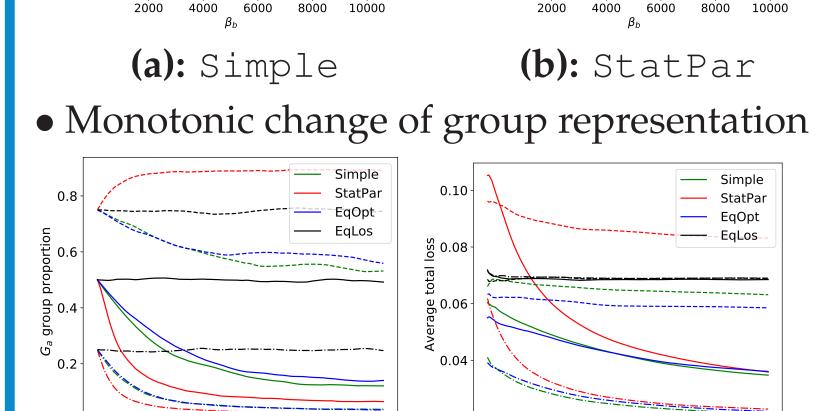




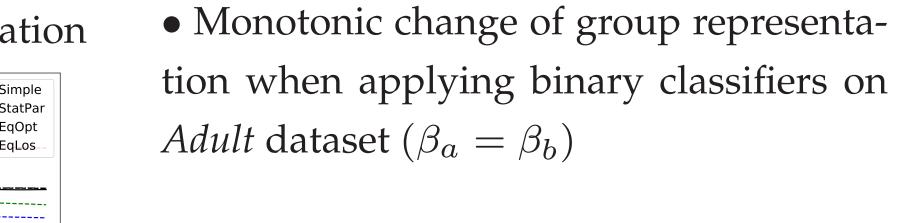
- Δ -fair set: all (θ_a, θ_b) s.t. $\left|\frac{\widetilde{N}_a}{\widetilde{N}_b} \frac{\beta_a}{\beta_b}\right| \leq \min\{\left|\frac{\widetilde{N}_a}{\widetilde{N}_b} \frac{\beta_a}{\beta_b}\right|\} + \Delta$
- Examples: $N_k(t+1) = N_k(t)\pi_k^2(\theta_k(t)) + \beta_k\pi_k^1(\theta_k(t))$ with $\Delta = \epsilon \frac{\beta_a}{\beta_b}$

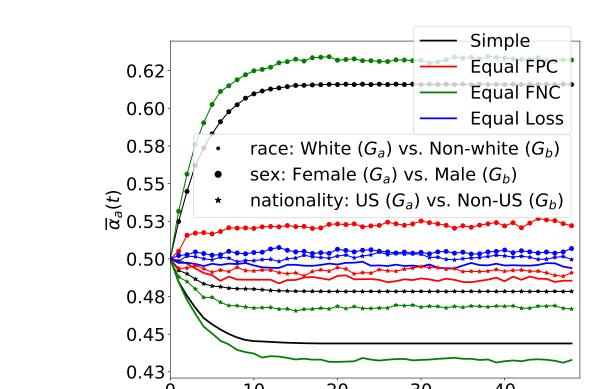
NUMERICAL RESULTS





• Effect of Δ -fair decisions ($\beta_a = \beta_b$)





CONCLUSIONS

- Group representation disparity can get exacerbated over time very easily under seemingly fair decisions.
- Exacerbation in representation disparity can accelerate when feature distributions are also reshaped by decisions.
- Develop a method of selecting a proper fairness criterion based on prior knowledge of participation dynamics.