

# A Sequentially Fair Mechanism for Multiple Sensitive Attributes

AAAI-24. Hu et al. [1]

Reviewer: Jihu Lee

IDEA lab  
Department of Statistics  
Seoul National University

April 15, 2025

- Using multi-marginal 2-Wasserstein barycenters, offers a closed-form solution of the learning problem (DP).
- Rewrite the optimal fair solution into a sequential form.

- Features:  $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^d$ , Target:  $Y \in \mathcal{Y} \subset \mathbb{R}$ .
- Sensitive:  $\mathbf{A} = (A_1, \dots, A_r) \in \mathcal{A}_1 \times \dots \times \mathcal{A}_r$  where  $\mathcal{A}_i = \{1, \dots, K_i\}$  with  $K_i \in \mathbb{N}$ .
- $A_{i:i+k} = (A_i, \dots, A_{i+k})$ .
- Predictor:  $f : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{Y}$ .  $\mathcal{F}$ : set of predictors  $f$ .
- $\nu_f$ : p.m. of  $f(\mathbf{X}, \mathbf{A})$ ,  $\nu_{f|\mathbf{a}}$ : p.m. of  $f(\mathbf{X}, \mathbf{A})|\mathbf{A} = \mathbf{a}$ .
- $F_{f|\mathbf{a}}(u) := \mathbb{P}(f(\mathbf{X}, \mathbf{A}) \leq u | \mathbf{A} = \mathbf{a})$ .
- $Q_{f|\mathbf{a}}(v) := \inf\{u \in \mathbb{R} : F_{f|\mathbf{a}}(u) \geq v\}$ .
- $\mathcal{R}(f) = \mathbb{E}(Y - f(\mathbf{X}, \mathbf{A}))^2$ .
- $f^*(\mathbf{X}, \mathbf{A}) = \mathbb{E}[Y | \mathbf{X}, \mathbf{A}]$

---

*p.m.=probability measure*

# Single Sensitive Case (previous result from [2])

- Sensitive  $A \in \{1, \dots, K\}$ ,  $p_a = \mathbb{P}(A = a)$ .
- **Fairness measure (Strong DP)**

$$\mathcal{U}(f) = \max_{a \in \mathcal{A}} \int_{u \in [0,1]} |Q_f(u) - Q_{f|a}(u)| du$$

- **Optimal fair predictor**

$$f_B = \arg \min_{f \in \mathcal{F}} \{\mathcal{R}(f) : \mathcal{U}(f) = 0\}$$

# Single Sensitive Case (previous result from [2])

- Minimizing  $\mathcal{R}(f) = \text{minimizing } \mathbb{E}(f^*(\mathbf{X}, A) - f(\mathbf{X}, A))^2$ .

$$\min_{f: \mathcal{U}(f)=0} \mathbb{E}(f^*(\mathbf{X}, A) - f(\mathbf{X}, A))^2 = \min_{\nu} \sum_{a \in [K]} p_a \mathcal{W}_2^2(\nu_{f^*|a}, \nu)$$

- A map to the minimizer  $\mu_{\mathcal{A}} : \mathcal{V} \rightarrow \mathcal{V}$ :

$$\mu_{\mathcal{A}}(\nu_{f^*}) := \min_{\nu} \sum_{a \in [K]} p_a \mathcal{W}_2^2(\nu_{f^*|a}, \nu) \text{ (Wasserstein Barycenter)}$$

- Closed-form solution of  $f$ :

$$f_B(\mathbf{x}, a) = \left( \sum_{a' \in \mathcal{A}} p_{a'} Q_{f^*|a'} \right) \circ F_{f^*|a}(f^*(\mathbf{x}, a))$$

- **Fairness measure (Strong DP)**

$$\mathcal{U}_i(f) = \max_{a_i \in \mathcal{A}_i} \int_{u \in [0,1]} |Q_f(u) - Q_{f|a_i}(u)| du$$

$$\mathcal{U}_{i,\dots,i+k}(f) = \mathcal{U}_{i:i+k}(f) = \mathcal{U}_i(f) + \dots + \mathcal{U}_{i+k}(f)$$

- **Optimal fair predictor (marginal)**

$$f_{B_i} = \arg \min_{f \in \mathcal{F}} \{\mathcal{R}(f) : \mathcal{U}_i(f) = 0\}$$

$$f_{B_i}(\mathbf{x}, \mathbf{a}) = \left( \sum_{a'_i \in \mathcal{A}_i} p_{a'_i} Q_{f^*|a'_i} \right) \circ F_{f^*|a_i}(f^*(\mathbf{x}, \mathbf{a}))$$

## Sequentially fair mechanism

- Optimal fair predictor

$$f_B = \arg \min_{f \in \mathcal{F}} \{ \mathcal{R}(f) : \mathcal{U}_{1:r}(f) = 0 \}$$

- Under some assumptions, the optimal fair predictor is (Proposition 4,5):

$$\begin{aligned} f_B(\mathbf{x}, \mathbf{a}) &= (f_{B_1} \circ \cdots \circ f_{B_r})(\mathbf{x}, \mathbf{a}) \\ &= (f_{B_{\sigma(1)}} \circ \cdots \circ f_{B_{\sigma(r)}})(\mathbf{x}, \mathbf{a}) \end{aligned}$$

## Single sensitive case [2]

$$f_{B_1}^{\varepsilon_1} \in \arg \min_{f \in \mathcal{F}} \{ \mathcal{R}(f) : \mathcal{U}_1(f) \leq \varepsilon_1 \cdot \mathcal{U}_1(f^*) \}$$

$$f_{B_1}^{\varepsilon_1}(\mathbf{X}, \mathbf{A}) := (1 - \varepsilon_1) \cdot f_{B_1}(\mathbf{X}, \mathbf{A}) + \varepsilon_1 \cdot f^*(\mathbf{X}, \mathbf{A})$$

## Multiple extension (This paper)

$$f_B^\varepsilon = \arg \min_{f \in \mathcal{F}} \left\{ \mathcal{R}(f) : \mathcal{U}(f) \leq \sum_{i=1, \dots, r} \varepsilon_i \cdot \mathcal{U}_i(f^*) \right\}$$

$$f_B^\varepsilon(\mathbf{X}, \mathbf{A}) = \left( f_{B_{\sigma(1)}}^{\varepsilon_{\sigma(1)}} \circ \dots \circ f_{B_{\sigma(r)}}^{\varepsilon_{\sigma(r)}} \right) (\mathbf{X}, \mathbf{A})$$



- $f^* \leftarrow \hat{f}$ : plug-in with any unfair ML model.

$$\widehat{f_{B_i}}(\mathbf{x}, \mathbf{a}) = \left( \sum_{a'_i \in \mathcal{A}_i} \hat{p}_{a'_i} \right) \circ \hat{Q}_{\hat{f}|a_i} \left( \hat{f}(\mathbf{x}, \mathbf{a}) \right)$$

$$\widehat{f_B} = \widehat{f_{B_1}} \circ \cdots \circ \widehat{f_{B_r}}$$

# Experiments

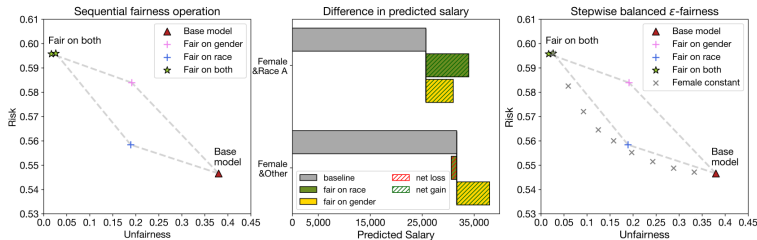


Figure 3: Applications on the *folktables* data set. Left pane, visualisation of the combined unfairness across two sensitive attributes and intermediate solutions rendering predictions fair on only one of them. Center pane, marginal changes to predicted income when rendering fair the predictions w.r.t. a single variable and the baseline predictions. Right pane, visualization of global metrics when correcting the score first for race, but keeping the average predicted salary of female individuals constant.

- [1] Hu, François, Philipp Ratz, and Arthur Charpentier. "A sequentially fair mechanism for multiple sensitive attributes." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 38. No. 11. 2024.
- [2] Chzhen, Evgenii, et al. "Fair regression with wasserstein barycenters." Advances in Neural Information Processing Systems 33 (2020): 7321-7331.

