



Algorithms and Analysis for Optimizing Robust Objectives in Fair Machine Learning

Cyrus Cousins
cbcousins@umass.edu

University of Amherst Massachusetts
Department of Computer Science

December 2023
New York, New York

Fairness in Operations and AI Workshop
Columbia University

What is Welfare-Centric Fair Machine Learning?

- Fair machine learning considers *multiple groups* ($\mathbf{x}_{1:g,1:m}, \mathbf{y}_{1:g,1:m}$)
- We can handle each group individually

◇ Empirical utility maximization

$$\hat{U}(h_\theta; \mathbf{x}_i, \mathbf{y}_i) \doteq \frac{1}{m} \sum_{j=1}^m U(h_\theta(\mathbf{x}_{i,j}), \mathbf{y}_{i,j}); \quad \forall i: \hat{\theta}_i \doteq \operatorname{argmax}_{\theta \in \Theta} U(h; \mathbf{x}_i, \mathbf{y}_i)$$

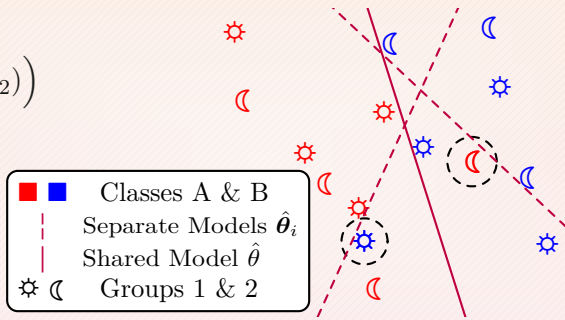
- What is the best classifier *overall*?

◇ Empirical welfare maximization

$$\hat{\theta} \doteq \operatorname{argmax}_{\theta \in \Theta} M(\hat{U}(h_\theta; \mathbf{x}_1, \mathbf{y}_1), \hat{U}(h_\theta; \mathbf{x}_2, \mathbf{y}_2))$$

- Welfare functions encode *social values*

- ◇ Optimize a *given welfare function* $M(\cdot)$
- ◇ Objectives specify tradeoffs!



Power Means and the Social Planner's Problem

- A *social planner* arranges society to the benefit of all
- How should we aggregate utility or disutility across groups?
- The power-mean for $p \in \mathbb{R}$ summarizes g values $s_{1:g}$ with *weights* $\mathbf{w}_{1:g}$ as

$$M_{p \neq 0}(\mathbf{s}; \mathbf{w}) \doteq \sqrt[p]{\sum_{i=1}^g \mathbf{w}_i s_i^p}$$

for $p \neq 0$, or

$$M_0(\mathbf{s}; \mathbf{w}) \doteq \exp \left(\sum_{i=1}^g \mathbf{w}_i \ln(s_i) \right)$$



- Fair welfare requires $p \leq 1$, extremes are interesting special cases
 - ◇ $p = 1$ is *weighted sum*, a.k.a. utilitarian welfare, over groups (well-studied case)
 - ◇ $p = 0$ is the *Nash social welfare* over groups
 - ◇ $p = -\infty$ limit is the *minimum* over groups (egalitarian or robust optimization)

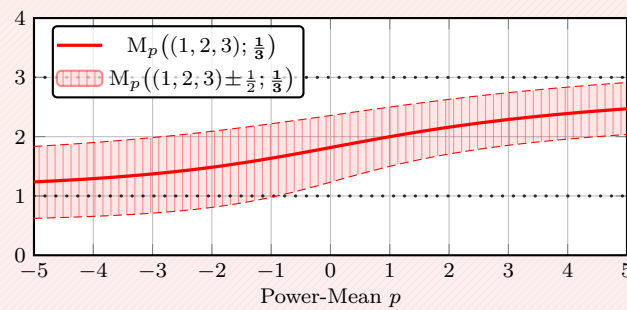
- Power-means are:

1. *Axiomatically Justified*
2. *Interpretable*

$$M_p(\mathbf{s}; \mathbf{w}) \text{ units match } s_{1:g}$$

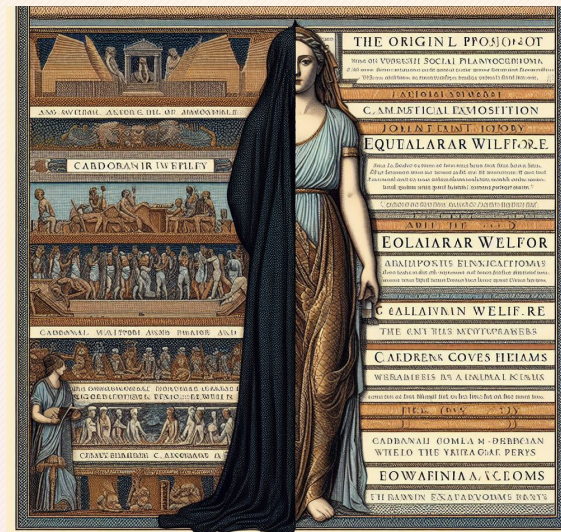
3. *Stochastically Stable*

(for $p \in [-\infty, 0) \cup [1, \infty]$)



John Rawls, the Original Position, and the Veil of Ignorance

- Justice, fairness, and societal wellbeing should be *objective concepts*
 - ◇ Should not depend on *our own identities*
- The “original position argument”
 - ◇ We should assess a situation from behind a “veil of ignorance”
- Rawls argues for *worst-case* robust or pessimistic analysis
 - ◇ Egalitarian welfare (malfare) is born!
 - ◇ Given utility $\mathbf{s} \in \mathbb{R}^g$, assess welfare as $M_{-\infty}(\mathbf{s}) = \min_{i=1}^g s_i$
 - ◇ Given disutility $\mathbf{s} \in \mathbb{R}^g$, assess malfare as $M_{\infty}(\mathbf{s}) = \max_{i=1}^g s_i$



This work: We pose the original problem setting as an *adversarial game*

- Robust fair objectives arise as solution concepts against specific adversaries

Egalitarianism and Adversarial Games

The Rawlsian Game

g : Number of inhabitants in the game world
 Θ : Parameter space of worlds Dæmon can create
 $\mathbf{s}(\theta) : \Theta \rightarrow \mathbb{R}_{0+}^g$: (Dis)utility vector of the inhabitants of world θ
 $\mathcal{S} \doteq \{\mathbf{s}(\theta) \mid \theta \in \Theta\} \subseteq \mathbb{R}_{0+}^g$: Set of *feasible* (dis)utility vectors the Dæmon can create
 $\mathcal{A}_{D\mathfrak{ae}} \doteq \mathcal{S}$: Dæmon's action space
 $\mathcal{A}_{Ang} \doteq \mathcal{W} \subseteq \Delta_g$: Angel's action space
 $P(\mathbf{s}, i) \doteq \pm(s_i, -s_i)$: 0-sum payoff function

Strategic gameplay (Dæmon goes first):

$$\arg \max_{\mathbf{s} \in \mathcal{A}_{D\mathfrak{ae}}} \min_{i \in \mathcal{A}_{Ang}} P_1(\mathbf{s}, i) = \arg \max_{\mathbf{s} \in \mathcal{S}} \min_{i \in 1, \dots, g} s_i$$



Weaker Adversaries, Robust Utilitarianism, and the Gini Social Welfare

The Constrained Rawlsian Game

$g, \Theta, \mathbf{s}(\theta), \mathcal{S}$: As above
 $\mathcal{A}_{D\mathfrak{ae}} \doteq \mathcal{S}$: Dæmon's action space
 $\mathcal{A}_{Ang} \doteq \mathcal{W} \subseteq \Delta_g$: Angel's action space (inhabitant i becomes *distribution* \mathbf{w})
 New 0-sum payoff function uses *expected* (dis)utility:

$$P(\mathbf{s}, \mathbf{w}) \doteq \langle \mathbf{w} \cdot \mathbf{s}, -\mathbf{w} \cdot \mathbf{s} \rangle$$

Strategic gameplay (Dæmon goes first):

$$\arg \max_{\mathbf{s} \in \mathcal{A}_{D\mathfrak{ae}}} \min_{\mathbf{w} \in \mathcal{A}_{Ang}} P_1(\mathbf{s}, \mathbf{w}) = \arg \max_{\mathbf{s} \in \mathcal{S}} \min_{\mathbf{w} \in \mathcal{W}} \mathbf{w} \cdot \mathbf{s}$$

A Few Special Cases

- Constant weights $\mathcal{W} = \{\mathbf{w}^*\}$: Optimize *utilitarian welfare* $M_1(\mathbf{s}; \mathbf{w}^*)$
- Norm ball $\mathcal{W} = \mathbf{w}^* + \{\mathbf{w} \mid \|\mathbf{w}\| \leq \gamma\}$: Optimize robust utilitarian welfare
- Weights *bounded from below* $\mathcal{W} \doteq \{\mathbf{w} \in \Delta_g \mid \mathbf{w} \succeq \gamma \mathbf{w}^*\}$: Optimize *utilitarian-maximin social welfare function* (UMSWF), i.e.,

$$\min_{\mathbf{w} \in \mathcal{W}} M_1(\mathbf{s}; \mathbf{w}) = \gamma M_1(\mathbf{s}; \mathbf{w}^*) + (1 - \gamma) M_{\mp\infty}(\mathbf{s})$$
- Given some *sorted weights vector* \mathbf{w}^\downarrow or \mathbf{w}^\uparrow , let \mathbf{s}^\uparrow denote *sorted* \mathbf{s}
 - ◇ $M_{\mathbf{w}^\downarrow}(\mathbf{s}) \doteq \mathbf{w}^\downarrow \cdot \mathbf{s}^\uparrow$ is *Gini social welfare*, and $M_{\mathbf{w}^\uparrow}(\mathbf{s}) \doteq \mathbf{w}^\uparrow \cdot \mathbf{s}^\uparrow$ is *Gini social malfare*
 - ◇ If Angel has actions $\mathcal{W} \doteq \{\pi(\mathbf{w}^\downarrow) \mid \pi \in \Pi_g\}$, where Π_g is all permutations $\{1, \dots, g\}$:

$$\max_{\mathbf{w} \in \mathcal{W}} M_1(\mathbf{s}; \mathbf{w}) = M_{\mathbf{w}^\downarrow}(\mathbf{s}) \quad \text{or} \quad \min_{\mathbf{w} \in \mathcal{W}} M_1(\mathbf{s}; \mathbf{w}) = M_{\mathbf{w}^\uparrow}(\mathbf{s})$$



A Divine Symmetry: Welfare and Malfare as Game-Theoretic Equilibria

A Game of Dæmonic Justice

- Suppose Dæmon wants to optimize *welfare* fairness concept $M_p(\mathbf{s}; \mathbf{w}^*)$
- Angel remains adversarial over some \mathcal{W}
- Payoff function

$$P(\mathbf{s}, \mathbf{w}) \doteq \langle M_p(\mathbf{s}; \mathbf{w}), -M_p(\mathbf{s}; \mathbf{w}) \rangle$$

- Strategic gameplay (Dæmon goes first):

$$\arg \max_{\mathbf{s} \in \mathcal{A}_{D\mathfrak{ae}}} \min_{\mathbf{w} \in \mathcal{A}_{Ang}} P_1(\mathbf{s}, \mathbf{w}) = \arg \max_{\mathbf{s} \in \mathcal{S}} \min_{\mathbf{w} \in \mathcal{W}} M_p(\mathbf{s}; \mathbf{w})$$

- ◇ Optimizes a *robust power-mean*!

A Game of Angelic Justice

- Suppose Angel wants to optimize $M_p(\mathbf{s}; \mathbf{w}^*)$ for some $p > 0$ with action space $\mathcal{W} = \Delta_g$
- We have the payoff function $P(\mathbf{s}, \mathbf{w}) \doteq \langle \mathbf{w} \cdot \mathbf{s}, M_p(\mathbf{s}; \mathbf{w}^*) \rangle$
- Angel strategy: Play $\mathbf{w}_i \propto \mathbf{w}_i^* s_i^{p-1}$
- Dæmon strategy: Select \mathbf{s} to optimize $\mathbf{s} \cdot \mathbf{w} = \sum_{i=1}^g \mathbf{w}_i^* s_i^p = M_p^p(\mathbf{s}; \mathbf{w}^*)$
- Play at this Nash equilibrium also optimizes a power-mean
- Angel can modify strategy to incorporate *robustness* $\mathbf{w}^* \in \mathcal{W}^*$



Fair Learning and Adversarial Optimization

We have reduced our fair and robust fair objectives to the common form

$$\arg \max_{\mathbf{s} \in \mathcal{A}_{D\mathfrak{ae}}} \min_{\mathbf{w} \in \mathcal{A}_{Ang}} P_1(\mathbf{s}; \mathbf{w}) = \arg \max_{\mathbf{s} \in \mathcal{S}} \min_{\mathbf{w} \in \mathcal{W}} M_p(\mathbf{s}; \mathbf{w})$$

Lemma 1 (Power-Mean Curvature)

For any $p \geq 1$, if $\mathbf{s} : \Theta \rightarrow \mathbb{R}^g$ is *convex*, then $M_p(\mathbf{s}(\theta); \mathbf{w}) : (\Theta \times \Delta_g) \rightarrow \mathbb{R}_{0+}$ is *convex* in θ and *concave* in \mathbf{w} .

For any $p \leq 1$, if $\mathbf{s} : \Theta \rightarrow \mathbb{R}^g$ is *concave*, then $M_p(\mathbf{s}(\theta); \mathbf{w}) : (\Theta \times \Delta_g) \rightarrow \mathbb{R}_{0+}$ is *concave* in θ and *convex* in \mathbf{w} .

Thus we can efficiently optimize robust fair objectives using *gradient ascent-descent* over $\mathbf{s} \in \mathcal{S}$

For indirect optimization tasks, we express the task in terms of parameter space Θ as

$$\arg \max_{\mathbf{s} \in \mathcal{A}_{D\mathfrak{ae}}} \min_{\mathbf{w} \in \mathcal{A}_{Ang}} P_1(\mathbf{s}; \mathbf{w}) = \arg \max_{\theta \in \Theta} \min_{\mathbf{w} \in \mathcal{W}} M_p(\mathbf{s}(\theta); \mathbf{w})$$

We can also efficiently optimize robust fair objectives using *gradient ascent-descent* over $\theta \in \Theta$

Robust Fair Objectives in Fair Machine Learning

- For any training point x and training label y , the *logistic loss* of model θ is

$$\ell(h_\theta(x), y) = \ln \left(\frac{1}{1 + \exp(-y \cdot x \cdot \theta)} \right)$$

- We then compute *disutility* for each group by averaging over training points

$$s_i(\theta) = \frac{1}{m_i} \sum_{j=1}^{m_i} \ell(h_\theta(\mathbf{x}_{i,j}), \mathbf{y}_{i,j})$$

- Observe: $s_i(\theta)$ is *strictly convex* over $\theta \in \Theta$: we can apply maximin methods
- Fair logistic regression objective is then

$$\operatorname{argmin}_{\theta \in \Theta} \sup_{\mathbf{w} \in \mathcal{W}} M_p \left(i \mapsto \frac{1}{m_i} \sum_{j=1}^{m_i} \ell(h_\theta(\mathbf{x}_{i,j}), \mathbf{y}_{i,j}); \mathbf{w} \right)$$

Robust Fair Objectives in Fair Allocation



A Most Curious Game

Combining welfare or malfare objectives and robustness yields the *Gini-power-mean family*

$$M_{\mathbf{w}^\downarrow, p}(\mathbf{s}) \doteq M_p(\mathbf{s}^\uparrow; \mathbf{w}^\downarrow) = \sqrt[p]{\sum_{i=1}^g \mathbf{w}_i^\downarrow (s_i^\uparrow)^p} \text{ for utility, } \quad \text{or}$$

$$M_{\mathbf{w}^\uparrow, p}(\mathbf{s}) \doteq M_p(\mathbf{s}^\uparrow; \mathbf{w}^\uparrow) = \sqrt[p]{\sum_{i=1}^g \mathbf{w}_i^\uparrow (s_i^\uparrow)^p} \text{ for disutility}$$

- Generalizes power-mean and Gini families

- ◇ Gini arises for $p = 1$

- ◇ Power-mean (unweighted) for $\mathbf{w}^\uparrow = \frac{1}{g} \mathbf{1}$ or $\mathbf{w}^\downarrow = \frac{1}{g} \mathbf{1}$

- Arises from power-mean axioms and a robust original position game!