

## Welfare, Malfare, and the Power Mean

The power-mean for  $p \in \mathbb{R} \setminus \{0\}$  summarizes  $g$  values  $\mathcal{S}_{1:g}$  with weights  $\mathbf{w}$ :

$$M_p(\mathcal{S}; \mathbf{w}) \doteq \sqrt[p]{\sum_{i=1}^g w_i \mathcal{S}_i^p}$$

Fair welfare:  $p \leq 1$ ,  $p = \infty$  is *minimum* over groups (egalitarian optimization)  
 Measure *overall wellbeing* given *utility values* (accuracy, income)

Fair malfare:  $p \geq 1$ ,  $p = \infty$  is *minimum* over groups (robust minimization)  
 Measure *overall illbeing* given *disutility values* (loss, harm)

Power-means are:

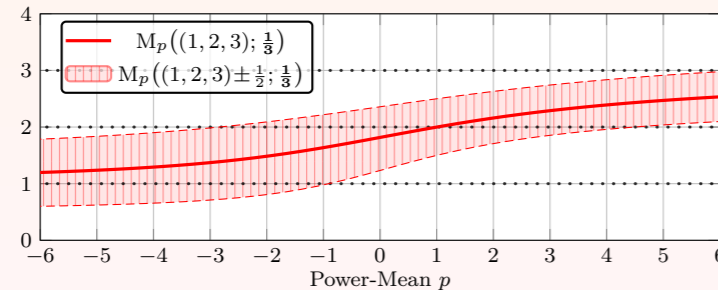
1. *Axiomatically Justified*

2. *Interpretable*

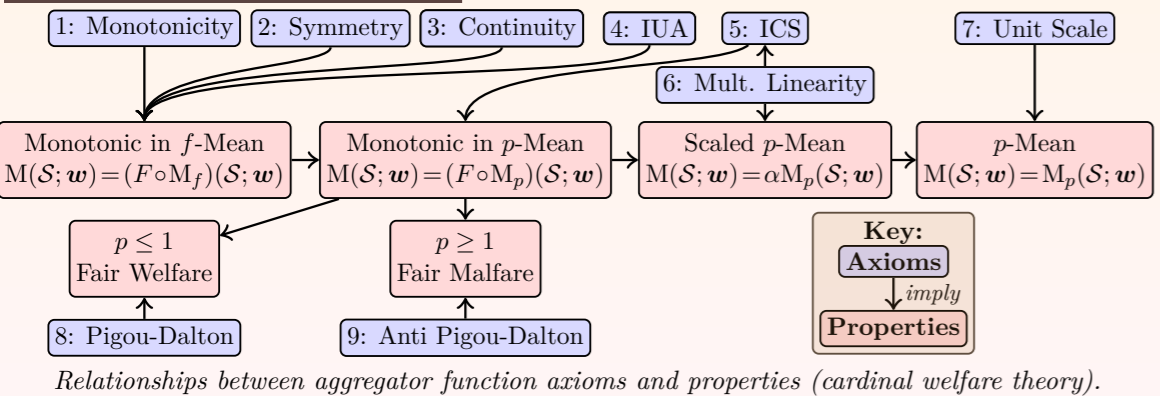
$M_p(\mathcal{S}; \mathbf{w})$  units match  $\mathcal{S}_{1:g}$

3. *Stochastically Stable*

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



## Axioms of Cardinal Welfare



## Estimating Malfare Values

1. Assuming only *monotonicity*:

Suppose  $\forall \omega \in \Omega: \hat{\mathcal{S}}(\omega) - \varepsilon(\omega) \leq \mathcal{S}(\omega) \leq \hat{\mathcal{S}}(\omega) + \varepsilon(\omega)$ . Then

$$M_p(\mathbf{0} \vee (\hat{\mathcal{S}} - \varepsilon); \mathbf{w}) \leq M_p(\mathcal{S}; \mathbf{w}) \leq M_p(\hat{\mathcal{S}} + \varepsilon; \mathbf{w}),$$

where  $\mathbf{a} \vee \mathbf{b}$  denotes the (elementwise) maximum.

2. Suppose range  $r$ . Then with probability at least  $1 - \delta$  over choice of  $\mathbf{x}$ :

$$|M_p(\mathcal{S}; \mathbf{w}) - M_p(\hat{\mathcal{S}}; \mathbf{w})| \leq r \sqrt{\frac{\ln \frac{2g}{\delta}}{2m}}$$

3. Suppose range  $r$  and variances  $\mathbb{V}_{\mathcal{D}_i}[\ell]$ . With probability at least  $1 - \delta$ :

$$|M_p(\mathcal{S}; \mathbf{w}) - M_p(\hat{\mathcal{S}}; \mathbf{w})| \leq \frac{r \ln \frac{2g}{\delta}}{3m} + \max_{i \in \{1, \dots, g\}} \sqrt{\frac{2 \mathbb{V}_{\mathcal{D}_i}[\ell] \ln \frac{2g}{\delta}}{m}}$$

N.b.: 2 & 3 hold for all fair malfare functions ( $p \geq 1$ ), but *not all* fair welfare functions.

## Empirical Malfare Minimization

Empirical risk and risk of hypothesis  $h$  given loss  $\ell$ :

$$\hat{R}(h; \ell, \mathbf{z}) \doteq \widehat{\mathbb{E}}_{(x,y) \in \mathbf{z}} [\ell(y, h(x))] \quad \& \quad R(h; \ell, \mathcal{D}) \doteq \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(y, h(x))]$$

We define *empirical malfare minimization* (EMM), given  $\mathbb{M}(\cdot; \mathbf{w})$ ,  $\mathcal{D}_{1:g}$ , and  $\mathbf{z}_{1:g}$ , with proxy and optimal models

$$\hat{h} \doteq \operatorname{argmin}_{h \in \mathcal{H}} \mathbb{M}(i \mapsto \hat{R}(h; \ell, \mathbf{z}_i); \mathbf{w}) \quad \& \quad h^* \doteq \operatorname{argmin}_{h \in \mathcal{H}} \mathbb{M}(i \mapsto R(h; \ell, \mathcal{D}_i); \mathbf{w})$$

Under what conditions is  $\hat{h}$  a *good proxy* for  $h^*$ ?

**Theorem 1** (Uniform Convergence of Malfare)

Suppose fair malfare  $\mathbb{M}_p(\cdot; \cdot)$  (i.e.,  $p \geq 1$ ), probability vector  $\mathbf{w} \in \mathbb{R}_+^g$ , loss function  $\ell: (\mathcal{Y} \times \mathcal{Y}) \rightarrow [0, r]$ , samples  $\mathbf{z}_i \sim \mathcal{D}_i^m$ , and hypothesis class  $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}$ .

Then with probability at least  $1 - \delta$  over choice of  $\mathbf{z}$ :

$$\sup_{h \in \mathcal{H}} |\mathbb{M}_p(i \mapsto R(h; \ell, \mathcal{D}_i); \mathbf{w}) - \mathbb{M}_p(i \mapsto \hat{R}(h; \ell, \mathbf{z}_i); \mathbf{w})| \leq \mathbb{M}_p\left(i \mapsto 2\hat{\mathbf{R}}_m(\ell \circ \mathcal{H}, \mathbf{z}_i) + 3r \sqrt{\frac{\ln \frac{g}{\delta}}{2m}}; \mathbf{w}\right)$$

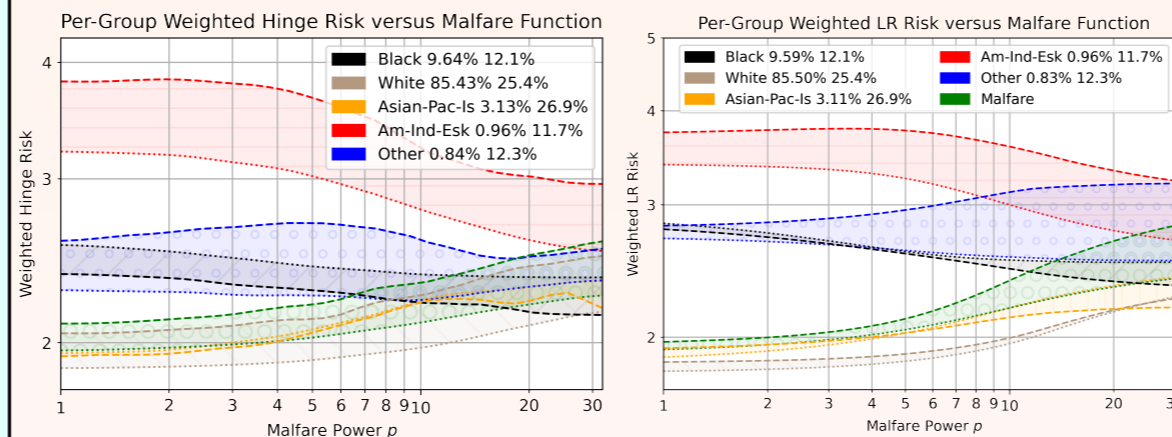
## Experiments

Training *linear models* on *adult* (census data) dataset

- Support vector machine (hinge loss)
- Logistic regression (cross entropy loss)
- Losses weighted by group-conditional label frequencies

Predict whether income exceeds 50,000\$ per annum

Minimize power-mean malfare over  $g = 5$  ethnic groups



Higher  $p \implies$  fairer model, closer to *egalitarianism*

$p = 1$  favors *large groups* (at the expense of minorities)

This is the *default* (if minority groups are even considered during training!)

Dire need for fairness-sensitive learning objectives

## Fair PAC Learnability

**Definition 2** (Fair-PAC Learnability)

Hypothesis class sequence  $\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \dots \subseteq \mathcal{X} \rightarrow \mathcal{Y}$  is fair-PAC-learnable w.r.t. loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{0+}$  if there exists a (randomized) algorithm  $\mathcal{A}$ , such that for all:

- sequence indices  $d$ ;
- $g$  instance distributions  $\mathcal{D}_{1:g}$ ;
- probability vectors  $\mathbf{w} \in \mathbb{R}_+^g$ ;
- malfares  $\mathbb{M}$  satisfying axioms 1-7+9;
- additive appx. errors  $\varepsilon > 0$ ;  $\mathcal{E}$
- failure probabilities  $\delta \in (0, 1)$ ;

$\mathcal{A}$  can identify a hypothesis  $\hat{h} \in \mathcal{H}$ , i.e.,  $\hat{h} \leftarrow \mathcal{A}(\mathcal{D}_{1:g}, \mathbf{w}, \mathbb{M}, \varepsilon, \delta, d)$ , where

- finite sample complexity:**  $\mathcal{A}(\mathcal{D}_{1:g}, \mathbf{w}, \mathbb{M}, \varepsilon, \delta, d)$  consumes no more than  $m(\varepsilon, \delta, d, g): (\mathbb{R}_+ \times (0, 1) \times \mathbb{N} \times \mathbb{N}) \rightarrow \mathbb{N}$  samples;  $\mathcal{E}$
- correctness:** with probability at least  $1 - \delta$ ,  $\hat{h}$  obeys

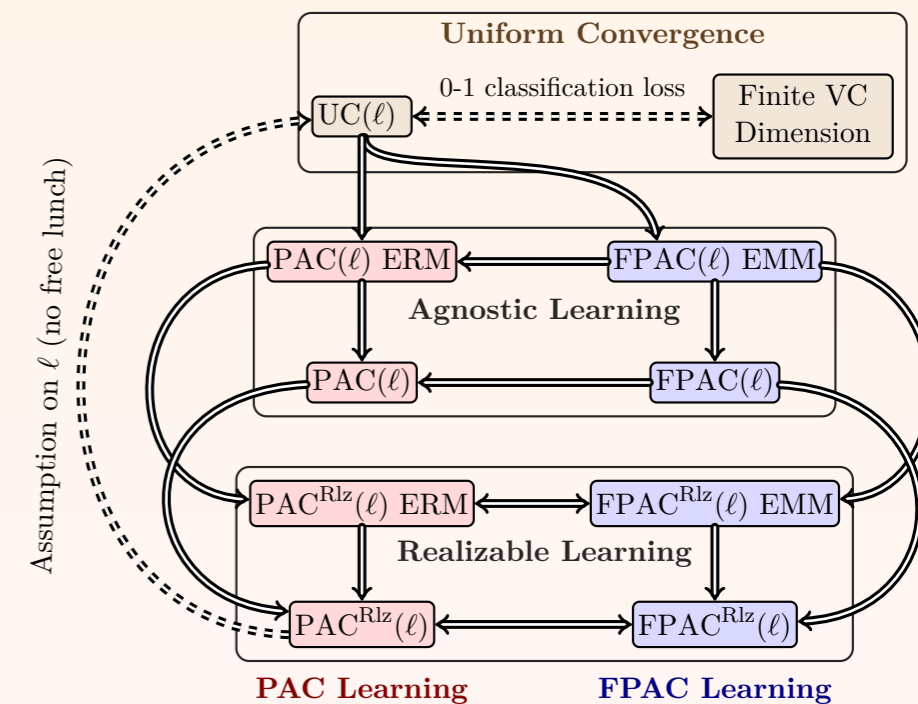
$$\mathbb{M}(i \mapsto R(\hat{h}; \ell, \mathcal{D}_i); \mathbf{w}) \leq \inf_{h^* \in \mathcal{H}} \mathbb{M}(i \mapsto R(h^*; \ell, \mathcal{D}_i); \mathbf{w}) + \varepsilon$$

The class of such fair-learning problems is denoted FPAC

If  $\forall d \in \mathbb{N}$ , the space of  $\mathcal{D}_{1:g}$  is restricted s.t.  $\inf_{h \in \mathcal{H}_d} \max_{i \in \{1, \dots, g\}} R(h; \ell, \mathcal{D}_i) = 0$ ,

then  $(\mathcal{H}, \ell)$  is *realizable-FPAC-learnable*, denoted  $(\mathcal{H}, \ell) \in \text{FPAC}^{\text{Rlz}}$

## Fundamental Theorem of Fair Statistical Learning



Implications between membership in PAC and FPAC classes

For fixed  $\ell$ ,  $\implies$  denotes *implication of membership* of some  $\mathcal{H}$  (i.e., containment)

Dashed implication arrows hold conditionally on  $\ell$

When the no-free-lunch assumption on  $\ell$  holds, the hierarchy collapses

Under realizability, some classes are known to coincide