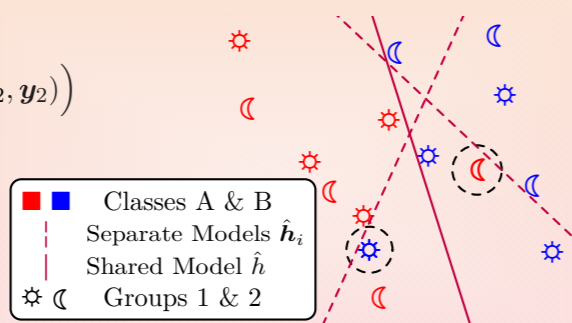


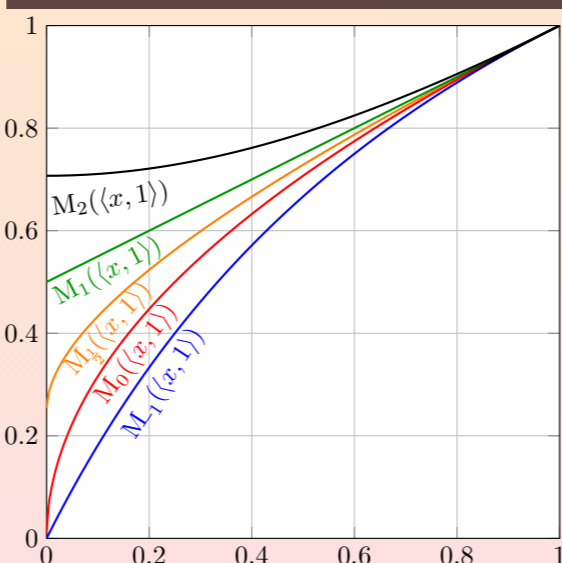


## What is Welfare-Centric Fair Machine Learning?

- 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
$$\hat{U}(h; \mathbf{x}_i, \mathbf{y}_i) \doteq \frac{1}{m} \sum_{j=1}^m U(h(\mathbf{x}_{i,j}), \mathbf{y}_{i,j}); \quad \forall i: \hat{h}_i \doteq \operatorname{argmax}_{h \in \mathcal{H}} U(h; \mathbf{x}_i, \mathbf{y}_i)$$
- What is the best classifier *overall*?
  - Empirical welfare maximization
$$\hat{h} \doteq \operatorname{argmax}_{h \in \mathcal{H}} W(\hat{U}(h; \mathbf{x}_1, \mathbf{y}_1), \hat{U}(h; \mathbf{x}_2, \mathbf{y}_2))$$
- Welfare functions encode *social values*
  - Optimize a *given welfare function*
  - Objectives specify tradeoffs!



## Small-Scale Behavior of Power-Means

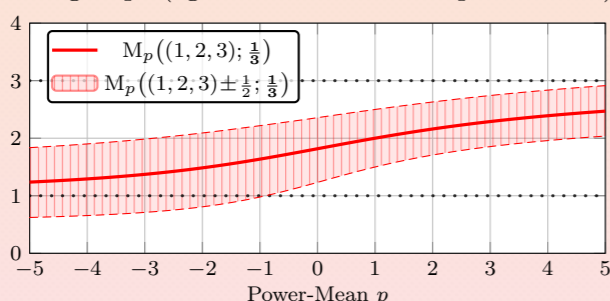


$$\frac{\partial}{\partial S_i} M_p(\mathcal{S}; \mathbf{w}) = w_i \underbrace{\left( \frac{S_i}{M_p(\mathcal{S}; \mathbf{w})} \right)^{p-1}}_{\text{Relative utility}}$$

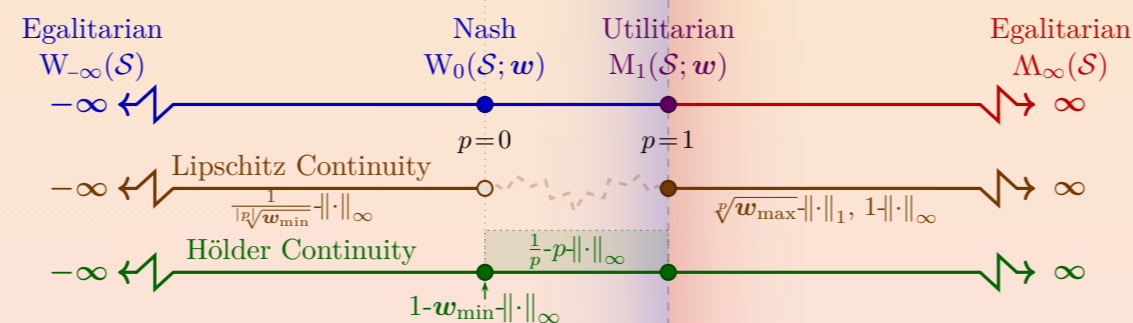
- Assuming *unit range*, power-means are:
- $p \geq 1$ :  $1 - \|\cdot\|_\infty$  Lipschitz
  - $p \in (0, 1)$ :  $\frac{1}{p} - p \|\cdot\|_\infty$  Hölder
  - $p = 0$ :  $1 - \mathbf{w}_{\min} \|\cdot\|_\infty$  Hölder
  - $p < 0$ :  $\frac{1}{\sqrt[p]{\mathbf{w}_{\min}}} \|\cdot\|_\infty$  Lipschitz
- The “difficult cases” occur as:
- $p \rightarrow 0$
  - $\mathbf{w}_{\min} \rightarrow 0$  for  $p < 1$

## Utilitarian, Egalitarian, and the Power-Mean Welfare

- The power-mean for  $p \in \mathbb{R}$  summarizes  $g$  values  $S_{1:g}$  with weights  $w_{1:g}$  as
 
$$M_{p \neq 0}(\mathcal{S}; \mathbf{w}) \doteq \sqrt[p]{\sum_{i=1}^g w_i S_i^p}, \quad M_0(\mathcal{S}; \mathbf{w}) \doteq \exp\left(\sum_{i=1}^g w_i \log(S_i)\right) = \prod_{i=1}^g S_i^{w_i}$$
- 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)



## Axiomatic Characterization of Welfare Classes



- Axiomatic characterization of welfare functions
  - Uniquely satisfied by  $W_p(\cdot; \mathbf{w})$  for  $p \leq 1$
- Additional axioms further restrict  $p$ 
  - Continuity properties vary by region
  - Impact on sample efficiency of learning

## A Hierarchy of Continuity Concepts

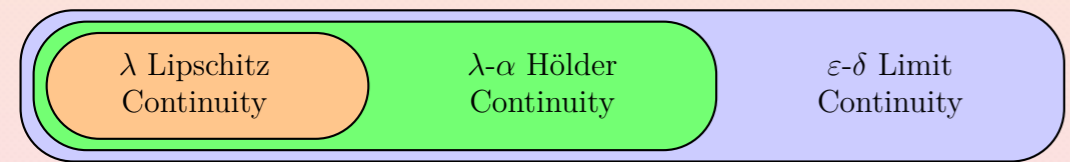
**Definition 1** (Hölder Continuity)  
 $M(\mathcal{S}; \mathbf{w})$  is Hölder continuous in  $\mathcal{S}$  with respect to norm  $\|\cdot\|_M$  if there exist some

- scale  $\lambda \geq 0$ ,
- power  $\alpha \in (0, 1]$ ,

such that for all  $\mathcal{S}, \mathcal{S}'$ , it holds that

$$|M(\mathcal{S}; \mathbf{w}) - M(\mathcal{S}'; \mathbf{w})| \leq \lambda \|\mathcal{S} - \mathcal{S}'\|_M^\alpha.$$

- Such a function is  $\lambda - \alpha \|\cdot\|_M$  Hölder continuous.
  - Bound the impact of *small changes*
- If  $\alpha = 1$ , it is  $\lambda \|\cdot\|_M$  Lipschitz continuous.
  - Bound the impact of *infinitesimal changes*



## Fair-PAC Learning

**Definition 2** (Fair-PAC Learning)  
 Suppose

- hypothesis class  $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}'$
- utility function  $U: \mathcal{Y}' \times \mathcal{Y} \rightarrow \mathbb{R}_{0+}$
- welfare class  $\mathcal{W} \subseteq \mathbb{R}_{0+}^g \rightarrow \mathbb{R}_{0+}$

$\mathcal{H}$  is fair-PAC-learnable if there exists an algorithm  $\mathcal{A}$  such that for any

- distributions  $\mathcal{D}_{1:g}$  over  $(\mathcal{X} \times \mathcal{Y})$
- welfare function  $W(\cdot; \mathbf{w}) \in \mathcal{W}$
- additive error  $\epsilon > 0$
- failure probability  $\delta \in (0, 1)$

$\mathcal{A}$  can identify a hypothesis  $\hat{h} \in \mathcal{H}$  such that

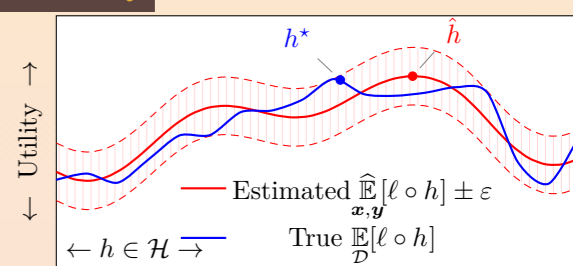
- $\mathcal{A}$  has  $m_{\mathcal{W}, \mathcal{H}}(\epsilon, \delta, W, g)$  sample complexity (per-group)
- with probability at least  $1 - \delta$ ,  $\hat{h}$  obeys

$$\underbrace{W\left(\mathbb{E}_{(x,y) \sim \mathcal{D}_1} [U(\hat{h}(x), y)], \dots; \mathbf{w}\right)}_{\text{Learned model welfare}} \geq \underbrace{\operatorname{argmax}_{h \in \mathcal{H}} W\left(\mathbb{E}_{(x,y) \sim \mathcal{D}_1} [U(h^*(x), y)], \dots; \mathbf{w}\right)}_{\text{Optimal model welfare}} - \epsilon$$

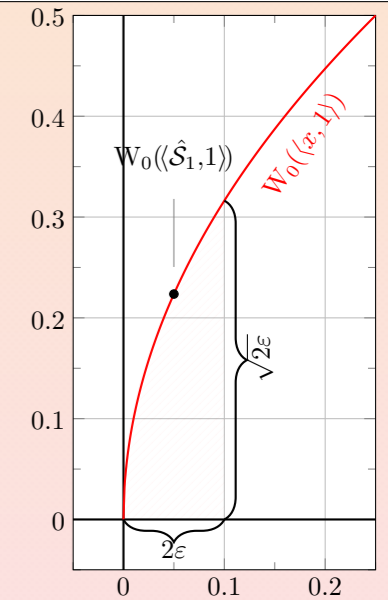
## Uniform Convergence and PAC-Learnability

- Suppose that for hypothesis class  $\mathcal{H}$ , with a sample  $(\mathbf{x}, \mathbf{y})$  of size  $m_{\mathcal{H}}(\epsilon, \delta)$ , it holds

$$\mathbb{P}_{\mathbf{x}, \mathbf{y}} \left( \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{\mathcal{D}} [\ell \circ h] - \widehat{\mathbb{E}}_{\mathbf{x}, \mathbf{y}} [\ell \circ h] \right| > \epsilon \right) < \delta$$



- Many ways to show this for various  $\mathcal{H}$ 
  - Vapnik-Chervonenkis dimension
  - Rademacher averages
- What does *uniform convergence* give us?
  - Asymptotic consistency of empirical utility maximizer  $\hat{h}$
  - Finite-sample convergence rate bounds
  - UC  $\implies$  PAC with sample complexity  $m_{\mathcal{H}}(\frac{\epsilon}{2}, \delta)$
- By  $\epsilon - \delta$  limit continuity alone:
  - Consistency of empirical welfare maximizer  $\hat{h}$
  - Finite-sample convergence rate bounds?
  - Convergence rate depends on welfare function
  - Asymptotic bounds in terms of  $\nabla_{\mathcal{S}} \cdot W(\mathcal{S}^*; \mathbf{w})$  for welfare-maximal utility vector  $\mathcal{S}^*$
- By how much does  $W(\cdot; \mathbf{w})$  magnify error?
  - Hölder continuity analysis



## Characterizing Fair-PAC Learnability

**Theorem 3** (Welfare Optimization Sample Complexity)  
 Suppose that for sample size  $m_{\mathcal{H}}(\epsilon, \delta)$ , it holds that

$$\mathbb{P}_{\mathbf{x}, \mathbf{y}} \left( \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{\mathcal{D}} [\ell \circ h] - \widehat{\mathbb{E}}_{\mathbf{x}, \mathbf{y}} [\ell \circ h] \right| > \epsilon \right) < \delta.$$

Then  $\mathcal{H}$  is FPAC-learnable with sample complexity  $m_{\mathcal{W}, \mathcal{H}}(\epsilon, \delta, W, g) \leq m_{\mathcal{H}}\left(\sqrt{\frac{\epsilon}{2\lambda}}, \frac{\delta}{g}\right)$ .

- Sample complexity of  $\epsilon - \delta$  learning bounded objectives is usually  $m_{\mathcal{H}}(\epsilon, \delta) \in \mathcal{O}\left(\frac{\ln \frac{1}{\delta}}{\epsilon^2}\right)$
- Fair-learning the class of *all weighted power-means* is thus
 
$$m_{\mathcal{W}, \mathcal{H}}(\epsilon, \delta, W, g) \leq m_{\mathcal{H}}\left(\sqrt{\frac{\epsilon}{2\lambda}}, \frac{\delta}{g}\right) \in \mathcal{O}\left(\frac{\lambda^\alpha \ln \frac{g}{\delta}}{\epsilon^\alpha}\right) \subseteq \underbrace{\mathcal{O}\left(\frac{\ln \frac{g}{\delta}}{\epsilon^{\frac{2}{\mathbf{w}_{\min}}}}\right)}_{\text{Worst-Case Bound!}}$$
- For any constant  $c \in (0, 1)$ , if  $\mathbf{w}_{\min} \geq \frac{c}{g}$  and  $|p| \geq c$ , then
 
$$m_{\mathcal{W}, \mathcal{H}}(\epsilon, \delta, W, g) \in \mathcal{O}\left(\frac{g^{\frac{2}{c}} \ln \frac{g}{\delta}}{\epsilon^{\frac{2}{c}}}\right) \subseteq \text{Poly}^{\frac{1}{c}}\left(\frac{1}{c}, \frac{1}{\epsilon}, \frac{1}{\delta}, \log \frac{1}{\delta}\right)$$

## A New Paradigm of Statistically Sound Fair Machine Learning

- Axiomatically characterize class of fair welfare functions
  - Fairness varies interpersonally, but “reasonable axioms” describe “reasonable people”
- The power-mean family (with  $p \leq 1$ ):  $M_p(\mathcal{S}; \mathbf{w}) \doteq \sqrt[p]{\sum_{i=1}^g w_i S_i^p}$
- Analyze continuity properties of fair welfare functions
  - Lipschitz and Hölder continuity:  $|M(\mathcal{S}; \mathbf{w}) - M(\mathcal{S}'; \mathbf{w})| \leq \lambda \|\mathcal{S} - \mathcal{S}'\|_M^\alpha$
- Fair-PAC learnability for all welfare functions  $W$  in class  $\mathcal{W}$ 
  - Uniform convergence  $\implies$  FPAC-Learnability