

To Pool or Not To Pool: Analyzing the Regularizing Effects of Group-Fair Training on Shared Models

Cyrus Cousins
cbcousins@umass.edu

I. Elizabeth Kumar
iekumar@brown.edu

Suresh Venkatasubramanian
suresh@brown.edu

Fairness and Overfitting

Given per-group samples $(\mathbf{x}, \mathbf{y}) = \mathbf{z}_i \in (\mathcal{X} \times \mathcal{Y})^{m_i}$ from \mathcal{D}^{m_i} , assume hypothesis class $\mathcal{H} \subseteq \mathcal{X} \rightarrow \mathcal{Y}$, loss function $\ell: \mathcal{Y}' \times \mathcal{Y} \rightarrow \mathbb{R}$

- Subgroups have potentially different distributions $\mathcal{D}_{1:g}$
- Worst-case generalization error can be analyzed on a per-group basis, but:
 - There may **too few samples** to compute bounds for minority groups
 - Data from large groups may **regularize overfitting** to small groups
- Our goal: Use **majority data** to bound **minority overfitting**

Setting: Malfare-based Fair Machine Learning

Empirical per-group risk:

$$\hat{R}_i(h, \mathbf{z}_i) \doteq \frac{1}{m_i} \sum_{j=1}^{m_i} \ell(h(x_j), y_j)$$

Choose a malfare function such as a power mean:

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

$p = 1$: \mathbf{w} -weighted risk minimization
 $p \rightarrow \infty$: *minimax fair learning*

Our problem is *empirical malfare minimization*:

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \Lambda(i \rightarrow \hat{R}_i(h); \mathbf{w})$$



Rademacher Averages

- Rademacher averages** $\mathfrak{R}_{m_i}(\ell \circ \mathcal{H}, \mathcal{D}_i)$ bound risk *generalization gap*

- Suppose range r loss
- Supremum Deviation (SD) Bound:** With probability at least $1 - \delta$:

$$\forall i \in 1, \dots, g: \sup_{h \in \mathcal{H}} |R_i(h) - \hat{R}_i(h)| \leq \varepsilon_i = 2\hat{\mathfrak{R}}_{m_i}(\ell \circ \mathcal{H}, \mathbf{z}_i) + r \sqrt{\frac{\ln \frac{g}{\delta}}{2m_i}}$$

- Can generalize this result to power-mean malfare

$$\sup_{h \in \mathcal{H}} \left| \Lambda_p(i \rightarrow R_i(h); \mathbf{w}) - \Lambda_p(i \rightarrow \hat{R}_i(h); \mathbf{w}) \right| \leq \max_{i \in 1, \dots, g} \varepsilon_i$$

Theoretical Restricted Hypothesis Classes

Let $\varepsilon_i \doteq r \sqrt{\frac{\ln \frac{1}{\delta}}{2m_i}}$ and $\eta_i \doteq 2\hat{\mathfrak{R}}_{m_i}(\ell \circ \mathcal{H}, \mathcal{D}_i) + \varepsilon_i$
We (pessimistically) upper-bound the objective value (w.h.p.)

$$\inf_{h' \in \mathcal{H}} \Lambda(j \mapsto \hat{R}(h', \mathbf{z}_j); \mathbf{w}) \leq \inf_{h' \in \mathcal{H}} \Lambda \left(j \mapsto \begin{cases} j \neq i & \hat{R}(h', \mathbf{z}_j) \\ j = i & R(h', \mathcal{D}_i) + \varepsilon_i \end{cases} \right)$$

and (optimistically) lower-bound the empirical malfare of all $h \in \mathcal{H}$ (w.h.p.)

$$\Lambda(j \mapsto \hat{R}(h, \mathbf{z}_j); \mathbf{w}) \geq \Lambda \left(j \mapsto \begin{cases} j \neq i & \hat{R}(h, \mathbf{z}_j) \\ j = i & R(h, \mathcal{D}_i) - \eta_i \end{cases} \right)$$

Set $\mathcal{H}_i^* \doteq \{h \in \mathcal{H}\}$, where

$$\Lambda \left(j \mapsto \begin{cases} j \neq i & \hat{R}(h, \mathbf{z}_j) \\ j = i & R(h, \mathcal{D}_i) - \eta_i \end{cases} \right) \leq \inf_{h' \in \mathcal{H}} \Lambda \left(j \mapsto \begin{cases} j \neq i & \hat{R}(h', \mathbf{z}_j) \\ j = i & R(h', \mathcal{D}_i) + \varepsilon_i \end{cases} \right)$$

Theorem 1. Assume as above; the following then hold:

- With probability at least $1 - 2\delta$ over choice of \mathbf{z}_i , it holds that $\hat{h} \in \mathcal{H}_i^*$.
- With probability at least $1 - 4\delta$ over choice of \mathbf{z}_i ,

$$|R(\hat{h}, \mathcal{D}_i) - \hat{R}(\hat{h}, \mathbf{z}_i)| \leq 2\hat{\mathfrak{R}}_{m_i}(\ell \circ \mathcal{H}_i^*, \mathcal{D}_i) + \varepsilon_i.$$

Empirical Restricted Hypothesis Classes

Take $\hat{\eta}_i \doteq 2\hat{\mathfrak{R}}_{m_i}(\ell \circ \mathcal{H}, \mathbf{z}_i) + 2\varepsilon_i$

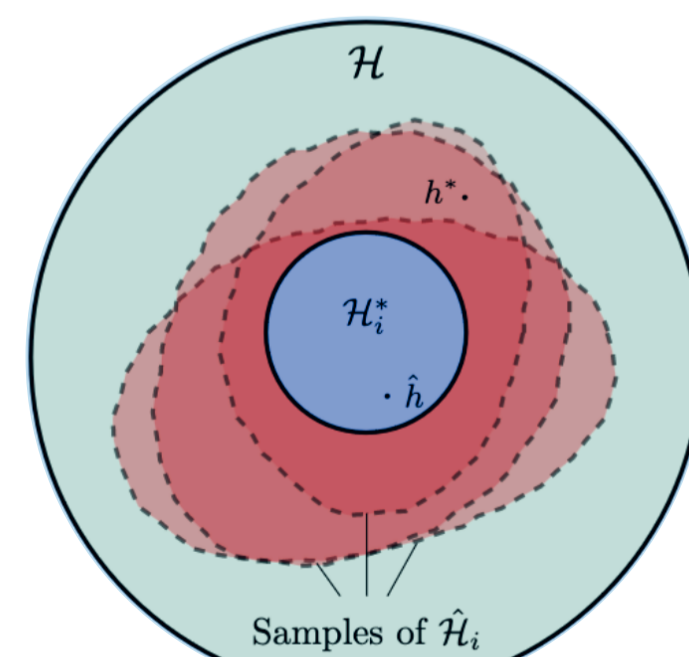
Set $\hat{\mathcal{H}}_i \doteq \{h \in \mathcal{H}\}$, where

$$\Lambda \left(j \mapsto \begin{cases} j \neq i & \hat{R}(h, \mathbf{z}_j) \\ j = i & \hat{R}(h, \mathbf{z}_i) - \hat{\eta}_i \end{cases} \right) \leq \inf_{h' \in \mathcal{H}} \Lambda \left(j \mapsto \begin{cases} j \neq i & \hat{R}(h', \mathbf{z}_j) \\ j = i & \hat{R}(h', \mathbf{z}_i) + 2\varepsilon_i \end{cases} \right)$$

Theorem 1. Assume as above; the following then hold:

- With probability at least $1 - 4\delta$ over choice of \mathbf{z}_i , it holds that $\hat{h} \in \mathcal{H}_i^* \subseteq \hat{\mathcal{H}}_i$.
- With probability at least $1 - 6\delta$, it holds that

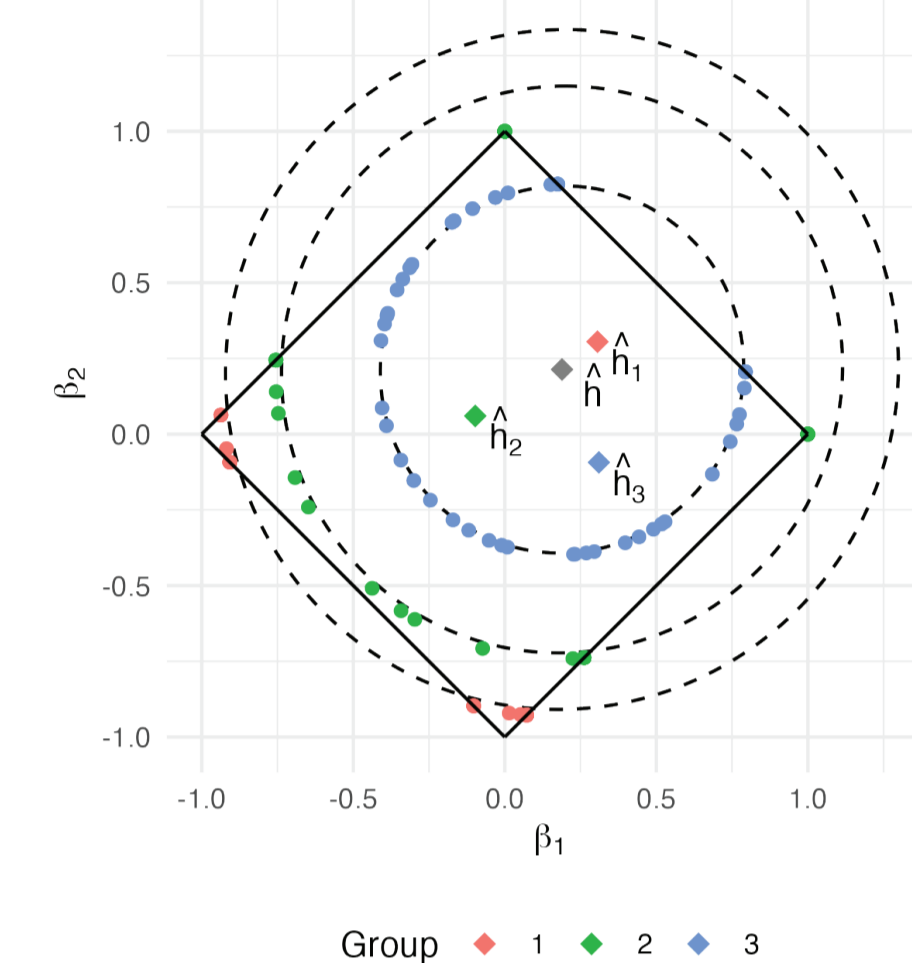
$$|R(\hat{h}, \mathcal{D}_i) - \hat{R}(\hat{h}, \mathbf{z}_i)| \leq 2\hat{\mathfrak{R}}_{m_i}(\ell \circ \hat{\mathcal{H}}_i, \mathbf{z}_i) + 2\varepsilon_i.$$



Visualization of unrestricted class \mathcal{H} , theoretical restricted class \mathcal{H}_i^* , and samples of empirical restricted class $\hat{\mathcal{H}}_i$ (varying \mathbf{z}_i).

Example with Linear Regression

- Take linear hypothesis class $\mathcal{B} \doteq \{\beta \in \mathbb{R}^2 : \|\beta\|_1 \leq 1\}$
- Sample (x_i, y_i) as $x_i \sim \text{Unif}([-1, 1]^2)$
 $y_i = x_i \cdot \beta_i + \text{Unif}([-1, 1])$
- Plot values of β which realize each supremum in the empirical Rademacher average for some Rademacher sample σ_k
- Points lie on either (the corner of) the ℓ_1 constraint boundary of \mathcal{B} or the restricted hypothesis constraint boundary of $\hat{\mathcal{H}}_i$



Experiments with Logistic Regression

$$\mathcal{X} = [-1, 1]^{15}$$

$$\mathcal{Y} = \pm 1$$

$$\mathcal{B} \doteq \{\beta \in \mathbb{R}^{15} : \|\beta\|_1 \leq 15\}$$

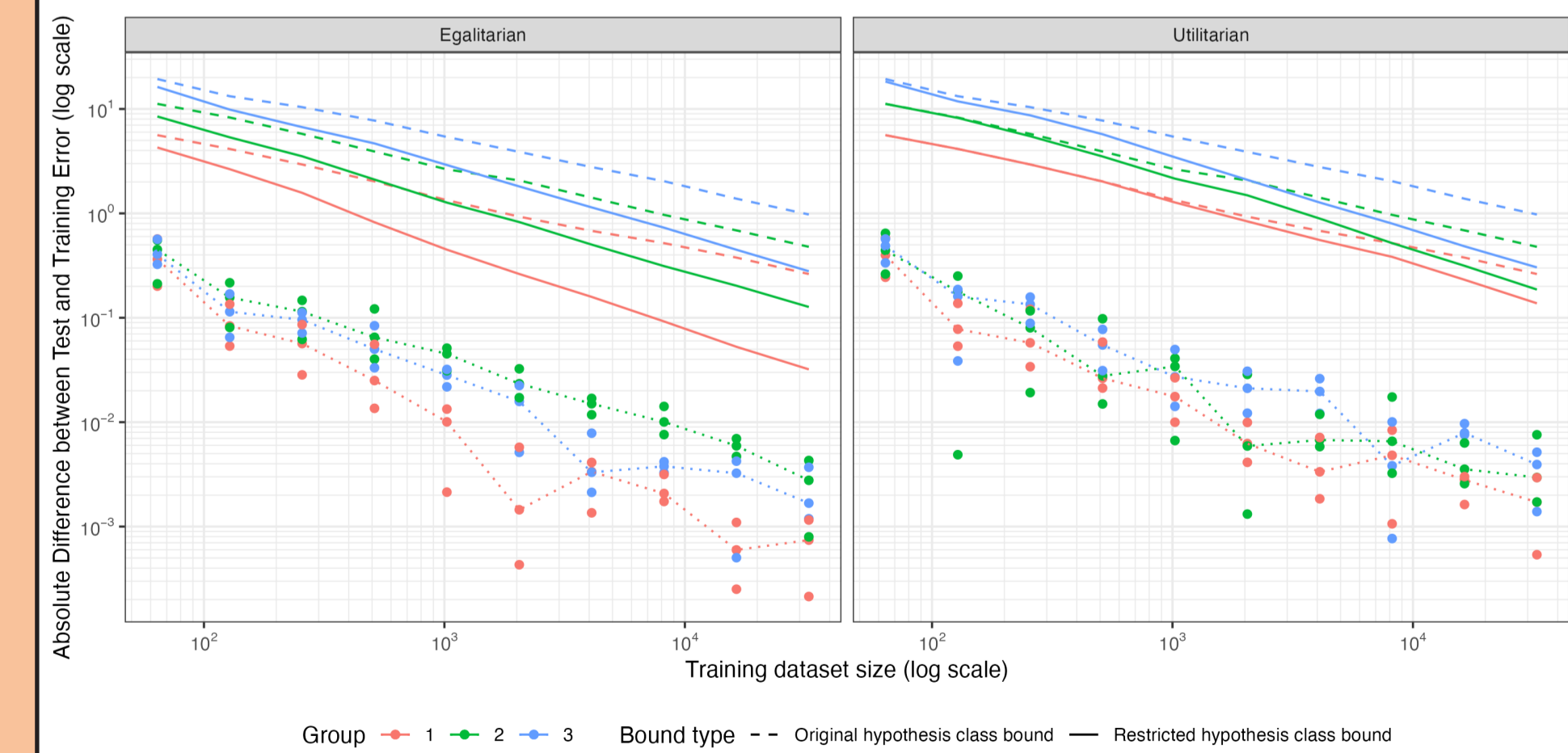
$$x \sim \text{Unif}(\mathcal{X})$$

$$\mathbb{P}(y = 1) = \text{logistic}(x \cdot \beta_i + \xi)$$

$$\hat{h} \doteq \operatorname{argmin}_{h \in \mathcal{H}} \Lambda(i \mapsto \hat{R}(h, \mathbf{z}_i))$$

$$\hat{R}(h, \mathbf{z}_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} \ln(1 + \exp(y_{i,j} \cdot h(x_{i,j})))$$

	Data proportion	True parameters
Group 1	75%	$\beta_i = 0.3$
Group 2	20%	$\beta_i = 0.1$
Group 3	5%	$\beta_i = 0.2$



This research was made possible in part by the generous support of the **Ford Foundation** and the **MacArthur Foundation**. Cyrus Cousins also wishes to acknowledge the **Center for Data Science at the University of Massachusetts Amherst**, where part of this work was conducted under a postdoctoral fellowship.

