



Classifiers and Margin [1]

Feature space \mathcal{X} , sensitive attributes \mathcal{S} , labels \mathcal{Y} .

Decision function $h \in \mathcal{H} \subseteq \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$.

We classify $x \in \mathcal{X}$ as:

$$H(x) = \arg \max_{y \in \mathcal{Y}} h(x, y) \quad \text{for } x \in \mathcal{X} .$$

Confidence margin of h for label y on input x :

$$\rho(h, x, y) = h(x, y) - \max_{y' \neq y} h(x, y') .$$

Group Fairness (Example of Equality of Opportunity [2])

Fairness level of $h \in \mathcal{H}$, for $(y, k) \in \mathcal{Y} \times \mathcal{S}$, for “desirable” labels y :

$$F_{(y,k)}(h, D) = \mathbb{P}(H(X) = Y \mid Y = y, S = k) - \mathbb{P}(H(X) = Y \mid Y = y) .$$

(Equalized odds, accuracy parity, and demographic parity have similar expressions.)

$$\text{Average fairness level: } \text{Fair}(h, D) = \frac{1}{|\mathcal{Y} \times \mathcal{S}|} \sum_{(y,k) \in \mathcal{Y} \times \mathcal{S}} F_{(y,k)}(h, D) .$$

Summary

The difference of fairness between private and optimal models vanishes since:

1. Group fairness notions are pointwise Lipschitz.
2. Models learned by output perturbation or DP-SGD converge to non-private one at a rate $O(\sqrt{p}/n\epsilon)$.

Main Assumption: Lipschitz Margins

For $x, y \in \mathcal{X} \times \mathcal{Y}$, there exists $L_{x,y}$ such that for all $h, h' \in \mathcal{H}$

$$\|\rho(h, x, y) - \rho(h', x, y)\| \leq L_{x,y} \|h - h'\| .$$

(Some) Group Fairness Notions are Pointwise Lipschitz

For $h, h' \in \mathcal{H}$, and any event E : $\left| \mathbb{P}(H(X) = Y \mid E) - \mathbb{P}(H'(X) = Y \mid E) \right| \leq \mathbb{E} \left(\frac{L_{X,Y}}{|\rho(h, X, Y)|} \mid E \right) \|h - h'\|$.

$$\implies |F_{(y,k)}(h, D) - F_{(y,k)}(h', D)| \leq \chi_{(y,k)}(h) \cdot \|h - h'\| , \quad \text{with } \chi_{(y,k)}(h) = \mathbb{E} \left(\frac{L_{X,Y}}{|\rho(h, X, Y)|} \mid Y = y, S = k \right) + \mathbb{E} \left(\frac{L_{X,Y}}{|\rho(h, X, Y)|} \mid Y = y \right) .$$

Fairness Loss due to Privacy Vanishes in $O(\sqrt{p}/n\epsilon)$

$$|F_{(y,k)}(h^*, D) - F_{(y,k)}(h^{\text{priv}}, D)| \leq \chi_{(y,k)}(h^{\text{ref}}) \cdot O(\sqrt{p}/n\epsilon) , \quad \text{for } h^{\text{ref}} \in \{h^{\text{priv}}, h^*\} .$$

This guarantees that the fairness level of h^{priv} is close to the one of h^* , even when the latter is unknown.

Numerical Results on Logistic Regression

With $\epsilon = 1$, $\delta = \frac{1}{n^2}$ on celebA ($n = 182k$) and folktables ($n = 1,600k$).

Three variants of the bound, depending on knowledge of h, h' :

- knowing only theoretical bound on $\|h^{\text{priv}} - h^*\|$,
- - - knowing empirical value of $\|h^{\text{priv}} - h^*\|$,
- ⋯ knowing actual values of h^{priv} and h^* .

More generally

With $F_k(h, D) = C_k^0 + \sum_{k'=1}^K C_k^{k'} \mathbb{P}(H(X) = Y \mid D_{k'})$,
(e.g. for equalized odds, accuracy parity, demographic parity...)

$$\chi_k(h) = \sum_{k'=1}^K C_k^{k'} \mathbb{E} \left(\frac{L_{X,Y}}{|\rho(h, X, Y)|} \mid D_{k'} \right) .$$

Private Empirical Risk Minimization [3]

Assume **strongly-convex** loss. Release an (ϵ, δ) -DP value:

$$h^{\text{priv}} \approx h^* \in \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{(x,s,y) \in D} \ell(h(x), y) , \quad (\text{ERM})$$

► Output Perturbation [3]:

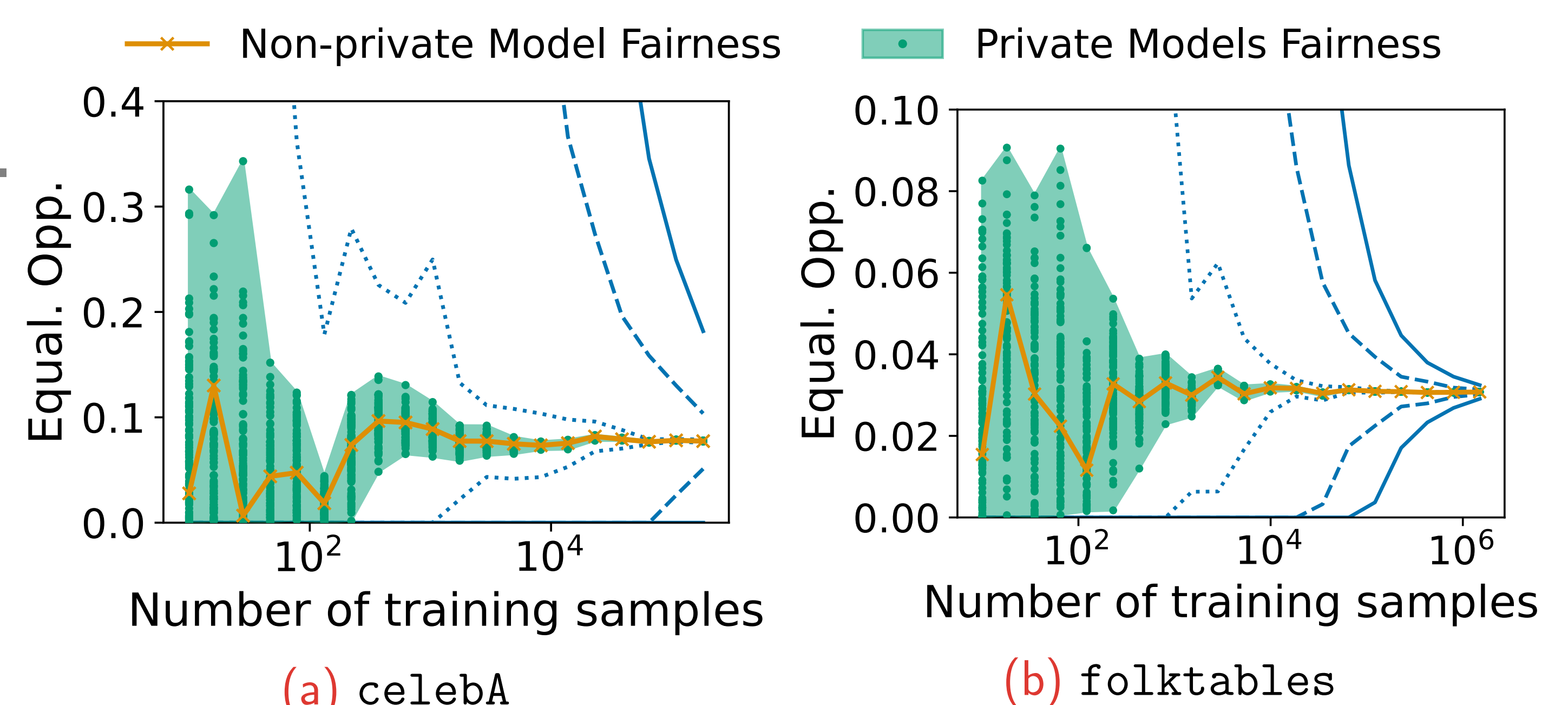
$$h^{\text{priv}} = h^* + \mathcal{N} \left(0, O \left(\frac{p}{n^2 \epsilon^2} \right) \right) .$$

► DP-SGD [4], compute for $t = 0 \dots, T - 1$:

$$h^{t+1} = h^t - \eta \left(\nabla f_t(h^t) + \mathcal{N} \left(0, O \left(\frac{pT}{n^2 \epsilon^2} \right) \right) \right) ,$$

and return $h^{\text{priv}} = h^T$.

In both cases: $\|h^{\text{priv}} - h^*\| = O \left(\frac{\sqrt{p}}{n\epsilon} \right)$ with high proba .



References

- [1] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. “Multi-Class Classification with Maximum Margin Multiple Kernel”. In: *ICML*. 2013.
- [2] Moritz Hardt et al. “Equality of Opportunity in Supervised Learning”. In: *NeurIPS*. 2016.
- [3] Kamalika Chaudhuri, Claire Monteleoni, and Anand D. Sarwate. “Differentially Private Empirical Risk Minimization”. In: *JMLR* (2011).
- [4] Raef Bassily, Adam Smith, and Abhradeep Thakurta. “Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds”. In: *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*. 2014.