Differential Privacy has Bounded Impact on Fairness in Classification

#### Paul Mangold, **Michaël Perrot**, Aurélien Bellet, Marc Tommasi Magnet Team

Workshop @ Comète on Ethical AI

November 23, 2023







イロト イポト イヨト イヨト



November 23, 2023 2 / 16

- 20

・ロト ・個ト ・モト ・モト



▲□▶ ▲圖▶ ▲理▶ ▲理▶ 二語:



◆□▶ ◆舂▶ ◆注≯ ◆注≯ ・注:



くロン 不得 とくほど 不良 とうしょう

M. Perrot



イロト イロト イヨト イヨト 三日



M. Perrot





















・ロト ・ 理 ト ・ ヨ ト ・ ヨ ・ うへぐ



(日)、(四)、(日)、(日)、

E

#### Notations

- Features space  $\mathcal{X}$ , labels space  $\mathcal{Y} = \{-1, 1\}$ , and sensitive attributes space  $\mathcal{S}$ .
  - Features: images, tabular data, graphs, ...
  - Labels: hired/not hired, profession, disease, ...
  - Sensitive attributes: gender, race, age, ...
- Decision function  $h: \mathcal{X} \mapsto \mathbb{R}$  taken in a set  $\mathcal{H}$ .
- Binary decision function  $H(x) = \begin{cases} 1 & \text{if } h(x) > 0, \\ -1 & \text{otherwise.} \end{cases}$
- A set  $D = \{(x_i, s_i, y_i)\}_{i=1}^n$  of *n* examples drawn i.i.d. a distribution  $\mathcal{D}$  over  $\mathcal{Z} = \mathcal{X} \times \mathcal{S} \times \mathcal{Y}$ .

**Goal:** Measure the fairness of the model *H*.

イロト イポト イヨト イヨト

Accuracy Parity [Zafar et al., 2017]

For all 
$$s \in S$$
,  $F_s(h, D) = \underbrace{\mathbf{P}(H(X) = Y | S = s)}_{\text{Accuracy of } h \text{ on the group } s} - \underbrace{\mathbf{P}(H(X) = Y)}_{\text{Accuracy of } h.}$ 





















イロト 不得下 不足下 不足下 一足

Accuracy Parity [Zafar et al., 2017]

3

イロト 不得下 不良下 不良下

Accuracy Parity [Zafar et al., 2017]

For all 
$$s \in S$$
,  $F_s(h, D) = \underbrace{\mathsf{P}(H(X) = Y | S = s)}_{\text{Accuracy of } h \text{ on the group } s} - \underbrace{\mathsf{P}(H(X) = Y)}_{\text{Accuracy of } h}$   

$$\begin{array}{c} & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & &$$

3

イロト 不得下 不良下 不良下



3

(日)、(四)、(日)、(日)、



3

(日)、(四)、(日)、(日)、



イロト イポト イヨト イヨト

3

M. Perrot



M. Perrot

3

## Differential Privacy [Dwork, 2006]

A randomized algorithm  $\mathcal{A}^{\text{priv}} : (\mathcal{X} \times \mathcal{S} \times \mathcal{Y})^n \to \mathcal{H}$  is  $(\epsilon, \delta)$ -private if for all neighboring datasets  $D, D' \in (\mathcal{X} \times \mathcal{S} \times \mathcal{Y})^n$  and all subsets of hypotheses  $\mathcal{H}' \subseteq \mathcal{H}$ 

 $\mathbb{P}(\mathcal{A}^{\mathsf{priv}}(D) \in \mathcal{H}') \leq \exp(\epsilon) \mathbb{P}(\mathcal{A}^{\mathsf{priv}}(D') \in \mathcal{H}') + \delta$ 



#### Gaussian Mechanism

Compute the  $\ell_2$ -sensitivity of the algorithm  $\mathcal{A}$ :

$$\Delta(\mathcal{A}) = \sup_{D \approx D'} \left\| \mathcal{A}(D) - \mathcal{A}(D') \right\|_{\mathcal{H}}$$



イロト 不得下 不同下 不同下

November 23, 2023 7 / 16

- 22

#### Gaussian Mechanism

Compute the  $\ell_2$ -sensitivity of the algorithm  $\mathcal{A}$ :

$$\Delta(\mathcal{A}) = \sup_{D pprox D'} \left\| \mathcal{A}(D) - \mathcal{A}(D') \right\|_{\mathcal{H}}$$

Add noise to  $\mathcal{A}(D)$ , calibrated to its sensitivity and the desired level of privacy:

$$\mathcal{A}^{\mathsf{priv}}(D) = \mathcal{A}(D) + \mathcal{N}\Big(0, rac{2\Delta(\mathcal{A})^2\log(1.25/\delta)}{\epsilon^2}\mathbb{I}_p\Big)$$



#### Gaussian Mechanism

Compute the  $\ell_2$ -sensitivity of the algorithm  $\mathcal{A}$ :

$$\Delta(\mathcal{A}) = \sup_{D pprox D'} \left\| \mathcal{A}(D) - \mathcal{A}(D') \right\|_{\mathcal{H}}$$

Add noise to  $\mathcal{A}(D)$ , calibrated to its sensitivity and the desired level of privacy:

$$\mathcal{A}^{\mathsf{priv}}(D) = \mathcal{A}(D) + \mathcal{N}\Big(0, rac{2\Delta(\mathcal{A})^2\log(1.25/\delta)}{\epsilon^2}\mathbb{I}_p\Big)$$



< 個 → < Ξ



イロト イポト イヨト イヨト

3



12

イロト イポト イヨト イヨト



M. Perrot

November 23, 2023 8 / 16

- 22

(日)、(四)、(日)、(日)、



$$\mathbf{P}(H(X) = Y | S = \text{red}) = \frac{4}{5}$$
$$\mathbf{P}(H(X) = Y) = \frac{9}{10}$$
$$F_{\text{red}}(h, D) = -\frac{1}{10}$$

November 23, 2023 9 / 16

- 22

イロト イポト イヨト イヨト



 $\mathbf{P}(H(X) = Y | S = \text{red}) = \frac{4}{5}$  $\mathbf{P}(H(X) = Y) = \frac{9}{10}$  $F_{\text{red}}(h, D) = -\frac{1}{10}$ 

What is the impact of differential privacy on fairness?

Pointwise Lipschitzness of Accuracy Parity

For all 
$$s \in S$$
,  $F_s(h, D) = \underbrace{\mathbf{P}(H(X) = Y | S = s)}_{\text{Accuracy of } h \text{ on the group } s} - \underbrace{\mathbf{P}(H(X) = Y)}_{\text{Accuracy of } h.}$ 

Let  $h, h' \in \mathcal{H}$  be two models

$$\left|F_{s}(h,D)-F_{s}(h',D)\right|\leq\chi_{s}(h,D)\left\|h-h'
ight\|_{\mathcal{H}}$$

イロト 不得下 イヨト イヨト

November 23, 2023

- 3

10/16

• 
$$\chi_s(h, D) = \mathbb{E}\left(\frac{L_X}{|h(X)|}\right) + \mathbb{E}\left(\frac{L_X}{|h(X)|} \mid S = s\right)$$

•  $L_X$  is a lipschitz constant of the model:  $|h(X) - h'(X)| \le L_X ||h - h'||_{\mathcal{H}}$ 

Pointwise Lipschitzness of Accuracy Parity

For all 
$$s \in S$$
,  $F_s(h, D) = \underbrace{\mathbf{P}(H(X) = Y | S = s)}_{\text{Accuracy of } h \text{ on the group } s} - \underbrace{\mathbf{P}(H(X) = Y)}_{\text{Accuracy of } h.}$ 

Let  $h, h' \in \mathcal{H}$  be two models

$$\left|F_{s}(h,D)-F_{s}(h',D)\right|\leq\chi_{s}(h,D)\left\|h-h'\right\|_{\mathcal{H}}$$

• 
$$\chi_s(h, D) = \mathbb{E}\left(\frac{L_X}{|h(X)|}\right) + \mathbb{E}\left(\frac{L_X}{|h(X)|} \mid S = s\right)$$

•  $L_X$  is a lipschitz constant of the model:  $|h(X) - h'(X)| \le L_X \|h - h'\|_{\mathcal{H}}$ 

**Key quantity:** Ratio between  $L_X(\downarrow)$  and the margin  $|h(X)|(\uparrow)$ .

イロト イポト イヨト イヨト

## Pointwise Lipschitzness of Accuracy Parity



3

(日)

# Output Perturbation $\mathcal{A}(D)$ is the following **optimization problem**

$$h^* = \operatorname*{arg\,min}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h; x_i, \underline{s}_i, y_i)$$

with  $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{S} \times \mathcal{Y} \to \mathbb{R}$  is  $\Lambda$ -lipschitz, and  $\mu$ -strongly-convex

$$\Delta(\mathcal{A}) = \frac{2\Lambda}{\mu n}$$

Using the Gaussian Mechanism, we can then release

$$h^{\mathsf{priv}} = h^* + \mathcal{N}\left(rac{8\Lambda^2\log(1.25/\delta)}{\mu^2 n^2 \epsilon^2}\mathbb{I}_p
ight)$$

A B + A B +
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

# Output Perturbation $\mathcal{A}(D)$ is the following **optimization problem**

$$h^* = \operatorname*{arg\,min}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h; x_i, \underline{s}_i, y_i)$$

with  $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{S} \times \mathcal{Y} \to \mathbb{R}$  is  $\Lambda$ -lipschitz, and  $\mu$ -strongly-convex

$$\Delta(\mathcal{A}) = \frac{2\Lambda}{\mu n}$$

Using the Gaussian Mechanism, we can then release

$$h^{\mathsf{priv}} = h^* + \mathcal{N}\left(rac{8\Lambda^2\log(1.25/\delta)}{\mu^2 n^2 \epsilon^2}\mathbb{I}_p
ight)$$

**Key result:** With probability at least  $1 - \zeta$ ,  $\|h^{\mathsf{priv}} - h^*\|_2 \leq \frac{\Lambda\sqrt{32p\log(1.25/\delta)\log(2/\zeta)}}{\mu n\epsilon}$ .

Bounded Distance

With probability at least 
$$1 - \zeta$$
,  $\|h^{\mathsf{priv}} - h^*\|_2 \leq \mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ .



990

イロト イロト イヨト イヨト 三日

#### Main result

- $|F_s(h, D) F_s(h', D)| \le \chi_s(h, D) ||h h'||_{\mathcal{H}}$
- With probability  $1 \zeta$ ,  $\left\| h^{\mathsf{priv}} h^* \right\|_2 \leq \mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$

With probability 
$$1 - \zeta$$
,  $\left| F_s(h^{\text{priv}}, D) - F_s(h^*, D) \right| \le \chi_s(h^{\text{priv}}, D) \mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ 

イロト 不得下 不足下 不足下 一足

#### Main result

- $|F_{s}(h,D) F_{s}(h',D)| \le \chi_{s}(h,D) ||h-h'||_{\mathcal{H}}$
- With probability  $1 \zeta$ ,  $\left\| h^{\mathsf{priv}} h^* \right\|_2 \leq \mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$

With probability 
$$1 - \zeta$$
,  $\left|F_s(h^{\text{priv}}, D) - F_s(h^*, D)\right| \le \chi_s(h^{\text{priv}}, D)\mathcal{O}\left(\frac{\sqrt{p}}{n\epsilon}\right)$ 

Auditing: The right hand side only depends on the private model!

イロト イポト イヨト イヨト

#### Conclusion

The story so far...

- Accuracy Parity is **pointwise lipschitz** and **margins** are key quantities.
- Using output perturbation the private and non-private models are close.
- Differential privacy has **bounded impact** on fairness.
- The squirrels are safe if they learned their model on sufficiently many examples.

A B A B A B
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

## Conclusion

The story so far...

- Accuracy Parity is **pointwise lipschitz** and **margins** are key quantities.
- Using output perturbation the private and non-private models are close.
- Differential privacy has **bounded impact** on fairness.
- The squirrels are safe if they learned their model on sufficiently many examples. but also...
  - Similar results for **other fairness measures** (e.g. Demographic Parity (with binary labels) [Calders et al., 2009], Equality of Opportunity [Hardt et al., 2016], Equalized Odds [Hardt et al., 2016]) and for **Accuracy**.
  - Multi-class, multi-groups problems, DP-SGD, tighter but harder to parse bounds.
  - A few **experiments** to check the tightness of our bounds.
  - A finite sample analysis showing that our results also hold in generalization.

3

イロト 不得下 イヨト イヨト



You can fetch the paper on arXiv!



Э

イロト イロト イヨト イヨ

M. Perrot

#### References I

All animal images were taken from commons.wikimedia.org. No animals were harmed during the preparation of this presentation.

- T. Calders, F. Kamiran, and M. Pechenizkiy. Building classifiers with independency constraints. In *2009 IEEE International Conference on Data Mining Workshops*, pages 13–18. IEEE, 2009.
- F. Ding, M. Hardt, J. Miller, and L. Schmidt. Retiring adult: New datasets for fair machine learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- C. Dwork. Differential privacy. In Encyclopedia of Cryptography and Security, 2006.
- M. Hardt, E. Price, and N. Srebro. Equality of opportunity in supervised learning. Advances in neural information processing systems, 29:3315–3323, 2016.
- M. B. Zafar, I. Valera, M. Gomez Rodriguez, and K. P. Gummadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th international conference on world wide web*, pages 1171–1180, 2017.

イロン 不得 とくせい 不良 とうせい

#### Experiments



Folktables [Ding et al., 2021],  $\ell_2$ -regularized logistic regression,  $(\epsilon, \delta) = (1, \frac{1}{n^2})$ 

November 23, 2023 18 / 16

#### Experiments



Folktables [Ding et al., 2021],  $\ell_2$ -regularized logistic regression,  $\delta = \frac{1}{n^2}$ , n = 1,498,050

November 23, 2023 19 / 16

< 17 ▶