38th Conference on Uncertainty in Artificial Intelligence Eindhoven, Netherlands August 1-5, 2022

uai2022

How unfair is private learning?

Amartya Sanyal, Yaxi Hu, Fanny Yang









Amartya

Yaxi



Fanny





Privacy and Fairness are both desirable properties in machine learning applications.









Privacy and Fairness are both desirable properties in machine learning applications.



Prior Work has mostly looked at the intersection:









Privacy and Fairness are both desirable properties in machine learning applications.



Prior Work has mostly looked at the intersection: Privacy and Accuracy: Kasiviswanathan et al. 2008, Feldman and Xiao 2014, Alon et. al., 2022.







Privacy and Fairness are both desirable properties in machine learning applications.





Prior Work has mostly looked at the intersection: Privacy and Accuracy: Kasiviswanathan et al. 2008, Feldman and Xiao 2014, Alon et. al., 2022. Fairness and Accuracy: Sagawa et. al. 2019, Du et al. 2021, Goel et. al. 2021.





Privacy and Fairness are both desirable properties in machine learning applications.





Prior Work has mostly looked at the intersection: Fairness and Accuracy: Sagawa et. al. 2019, Du et al. 2021, Goel et. al. 2021.





- Privacy and Accuracy: Kasiviswanathan et al. 2008, Feldman and Xiao 2014, Alon et. al., 2022.

THIS WORK: The interaction of Privacy and Fairness of nearly accurate algorithms.



Differential Privacy

Differential Privacy 0 0 0 0 0 0 0 0 1 1 1 1 1 1 7 Neighbor 3 3 3 3 3 3 3 3





Differential Privacy OOOOOO OOOOO A 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 Veighbor

 $(\epsilon, \delta) - \mathsf{DP}$ Algorithm



Differential Privacy 000000 2222222 333333 $(\epsilon, \delta) - \mathsf{DP}$ Algorithm



4





















Thrillers

Superhero











B&W	Mimes	Silent	Puppet	Oste
4%	4%	4%	4%	4%





B&W	Mimes	Silent	Puppet	Oste
4%	4%	4%	4%	4%



(Un) Fairness (Accuracy Discrepancy) ML Problem: Is the movie safe to watch for kids ?



B&W	Mimes	Silent	Puppet	Oste
4%	4%	4%	4%	4%



(Un) Fairness (Accuracy Discrepancy) ML Problem: Is the movie safe to watch for kids ?



B&W	Mimes	Silent	Puppet	Oste
4%	4%	4%	4%	4%
65%	75%	80%	80%	50



(Un) Fairness (Accuracy Discrepancy) ML Problem: Is the movie safe to watch for kids ?



B&W	Mimes	Silent	Puppet	Oste
4%	4%	4%	4%	4%
65%	75%	80%	80%	50
	Minorit	y Error =	70%	



(Un) Fairness (Accuracy Discrepancy) ML Problem: Is the movie safe to watch for kids?



Total Error = 18%

B&W	Mimes	Silent	Puppet	Oste	
4%	4%	4%	4%	4%	
65%	75%	80%	80%	50	
Minority Error = 70%					





(Un) Fairness (Accuracy Discrepancy) ML Problem: Is the movie safe to watch for kids?







Minority subpopulations

B&W	Mimes	Silent	Puppet	Oste
4%	4%	4%	4%	4%
65%	75%	80%	80%	50
	Minorit	y Error =	70%	

Accuracy Discrepancy = Minority Error - Total Error





(Un) Fairness (Accuracy Discrepancy) ML Problem: Is the movie safe to watch for kids?







Minority subpopulations

Minority Error = 70%					
65%	75%	80%	80%	50	
4%	4%	4%	4%	4%	
B&W	Mimes	Silent	Puppet	Oste	

Accuracy Discrepancy = 70 - 18 = 52%





40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.



40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.

• Subpopulation 1: Eyeglasses, bangs, ..., pointy nose.



40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.

- Subpopulation 1: Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,....,pointy noise.



40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.

- **Subpopulation 1**: Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,....,pointy noise.
- . . .
- . . .
- Subpopulation 2⁴⁰: No eyeglasses, no bangs,..., no pointy nose.


Example dataset CelebA

40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.

- Subpopulation 1: Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,....,pointy noise.
- . . .
- . . .
- Subpopulation 2⁴⁰: No eyeglasses, no bangs,..., no pointy nose.





Example dataset CelebA

40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.

- Subpopulation 1: Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,....,pointy noise.
- . . .
-
- Subpopulation 2⁴⁰: No eyeglasses, no bangs,..., no pointy nose.





Example dataset CelebA

40 binary attributes with each image



Eyeglasses

Bangs

Pointy Noise

40 binary attributes -> **2**⁴⁰ **subpopulations**.

- Subpopulation 1: Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,....,pointy noise.
- . . .
- . . .
- Subpopulation 2⁴⁰: No eyeglasses, no bangs,..., no pointy nose.





























Main Contribution: We prove this trend in a model-agnostic setting





Main Contribution: We prove this trend in a model-agnostic setting for long-tailed distribution.





Main Contribution:

• Error

• Error $\operatorname{err}(A, \Pi, F) =$

• Error $err(A, \Pi, F) =$ Learning Algorithm Data Distribution

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error err $(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ **Data Distribution**

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error err $(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ **Data Distribution**

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error err $(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m f \sim F h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ **Data Distribution**

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error err $(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m, f \sim F h \sim A(S_f)$ and $x \sim \Pi_{p,N}$ **Data Distribution**

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error err $(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ **Data Distribution**

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error $err(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m$, $f \sim F$, $h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ Data Distribution

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error $err(A, \Pi, F) = \square[h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ Data Distribution

Accuracy Discrepancy

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error $\operatorname{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$ Learning Algorithm Probability is over $S \sim \Pi^m$, $f \sim F$, $h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ **Data Distribution**

Accuracy Discrepancy

$\Gamma(A, \Pi, F) = \operatorname{err}_{\operatorname{Minority}}(A, \Pi, F) - \operatorname{err}(A, \Pi, F)$

Definitions of error and fairness Prior distribution over labelling functions $\subseteq Y^X$ • Error $\operatorname{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$ Learning Algorithm / Probability is over $S \sim \Pi^m$, $f \sim F$, $h \sim A(S_f)$, and $x \sim \Pi_{p,N}$ **Data Distribution**

Accuracy Discrepancy

Marginalised over minority subpopulations $\Gamma(A, \Pi, F) = \operatorname{err}_{\operatorname{Minority}}(A, \Pi, F) - \operatorname{err}(A, \Pi, F)$

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

(Privacy) Increases with privacy parameter ϵ .

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

(Privacy) Increases with privacy parameter ϵ .

N: # Minority subpopulations *m*: # Training points

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\xrightarrow{N} \rightarrow c$ as $N, m \rightarrow \infty$. M

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

• (Privacy) Increases with privacy parameter ϵ .

N: # Minority subpopulations *m*: # Training points

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\xrightarrow{N} \rightarrow c$ as $N, m \rightarrow \infty$. M

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .

N: # Minority subpopulations *m*: # Training points

• (Long-tailed) Increases with (relative) # of minority subpopulations c.

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\xrightarrow{N} \rightarrow c$ as $N, m \rightarrow \infty$. M

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .

N: # Minority subpopulations *m*: # Training points F: Label prior

• (Long-tailed) Increases with (relative) # of minority subpopulations c.

Privacy at the cost of fairness

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\xrightarrow{N} \rightarrow c$ as $N, m \rightarrow \infty$. (Label prior Entropy) Define $||F||_{\infty}$ =

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .

$$= \max_{x,y} \mathbb{P}_{f \sim F} \left[f(x) = y \right]$$

N: # Minority subpopulations *m*: # Training points F: Label prior

• (Long-tailed) Increases with (relative) # of minority subpopulations c.





Privacy at the cost of fairness

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\xrightarrow{N} \rightarrow c$ as $N, m \rightarrow \infty$. (Label prior Entropy) Define $||F||_{\infty}$ =

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .
- (Label prior) Increases with entropy of the label prior.

$$= \max_{x,y} \mathbb{P}_{f \sim F} \left[f(x) = y \right]$$

N: # Minority subpopulations *m*: # Training points F: Label prior

• (Long-tailed) Increases with (relative) # of minority subpopulations c.







 Most datasets follow longtailed distribution.

- Most datasets follow longtailed distribution.
 - With multiple small subpopulations.

- Most datasets follow longtailed distribution.
 - With multiple small subpopulations.
 - Individually: small mass.

- Most datasets follow longtailed distribution.
 - With multiple small subpopulations.
 - Individually: small mass.
 - Together: sizeable mass.

- Most datasets follow long-tailed distribution.
 - With multiple small subpopulations.
 - Individually: small mass.
 - Together: sizeable mass.



The number of examples by object class in the SUN dataset

- Most datasets follow longtailed distribution.
 - With multiple small subpopulations.
 - Individually: small mass.
 - Together: sizeable mass.



The number of examples by object class in the SUN dataset





An algorithm A is (ϵ, δ) -DP if

An algorithm A is $(\epsilon,\delta)\text{-}\mathsf{DP}$ if

$\mathbb{P}\left[A(S_1) \in Q\right] \le e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

An algorithm A is $(\epsilon,\delta)\text{-}\mathsf{DP}$ if

$\mathbb{P}\left[A(S_1) \in Q\right] \le e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space,

An algorithm A is $(\epsilon,\delta)\text{-}\mathsf{DP}$ if

$\mathbb{P}\left[A(S_1) \in Q\right] \le e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space, for all neighboring datasets S_1, S_2 .

An algorithm A is $(\epsilon,\delta)\text{-}\mathsf{DP}$ if

$\mathbb{P}\left[A(S_1) \in Q\right] \le e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space, for all neighboring datasets S_1, S_2 .

Neighboring datasets differ in only one entry.

An algorithm A is (ϵ, δ) -DP if

• Privacy parameter $\epsilon \in (0,\infty)$ $\mathbb{P}\left[A(S_1) \in Q\right] \leq \mathcal{E} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space, for all neighboring datasets S_1, S_2 .

Neighboring datasets differ in only one entry.

An algorithm A is (ϵ, δ) -DP if

$\mathbb{P}\left[A(S_1) \in Q\right] \leq e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space, for all neighboring datasets S_1, S_2 .

Neighboring datasets differ in only one entry.

• Privacy parameter $\epsilon \in (0,\infty)$

• Smaller is more private

An algorithm A is (ϵ, δ) -DP if

$\mathbb{P}\left[A(S_1) \in Q\right] \leq e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space, for all neighboring datasets S_1, S_2 .

Neighboring datasets differ in only one entry.

• Privacy parameter $\epsilon \in (0,\infty)$

Smaller is more private

Probability of catastrophic exposure



An algorithm A is (ϵ, δ) -DP if

$\mathbb{P}\left[A(S_1) \in Q\right] \leq e^{\epsilon} \mathbb{P}\left[A(S_2) \in Q\right] + \delta$

for all subsets Q of the output space, for all neighboring datasets S_1, S_2 .

Neighboring datasets differ in only one entry.

• Privacy parameter $\epsilon \in (0,\infty)$

• Smaller is more private

- Probability of catastrophic exposure
- Smaller is safer























Law School (Tabular) With Random Forest



Theorem B





Defining subpopulations



