

Learning to Walk Impartially on the Pareto Frontier of Fairness, Privacy, and Utility

Mohammad Yaghini, Patty Liu, Franziska Boenisch, Nicolas Papernot

University of Toronto & Vector Institute

Motivation: Specification Problem

- We want to deploy a model for chest X-rays to all regional hospitals
- Building the model is outsourced to a private company
- Regulators have privacy and fairness concerns
- They want to specify acceptable levels of trustworthiness guarantees for the model

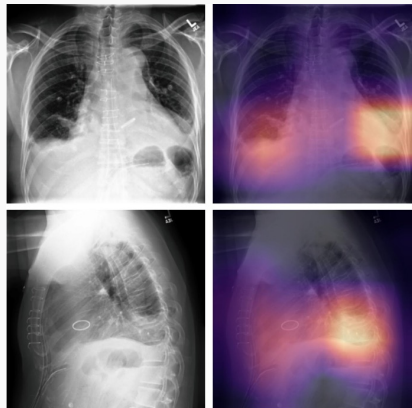


Figure 1: CheXpert

Motivation: Specification Problem

- We want to deploy a model for chest X-rays to all regional hospitals
- Building the model is outsourced to a private company
- Regulators have privacy and fairness concerns
- They want to specify acceptable levels of trustworthiness guarantees for the model

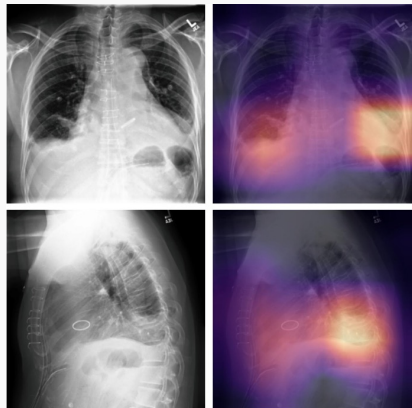


Figure 1: CheXpert

What are trustworthiness objectives for ML?

Values	Objective Examples	Mechanisms
Utility	Accuracy	Architecture search, optimizer search, etc.
Privacy	Differential Privacy (DP Loss) Unlearning	DP mechanisms: Noising, Randomized Response, etc.
Fairness	Demographic Parity (DemParity) Equality of Odds Disparate Impact	DemParity processors and regularizers

and many more (interpretability, robustness to distribution shifts/adversarial examples, etc.)

What are trustworthiness objectives for ML?

Values	Objective Examples	Mechanisms
Utility	Accuracy	Architecture search, optimizer search, etc.
Privacy	Differential Privacy (DP Loss) Unlearning	DP mechanisms: Noising, Randomized Response, etc.
Fairness	Demographic Parity (DemParity) Equality of Odds Disparate Impact	DemParity processors and regularizers

and many more (interpretability, robustness to distribution shifts/adversarial examples, etc.)

Privacy Objective: Differential Privacy

Definition $((\epsilon, \delta)$ -Differential Privacy)

Let $\mathcal{M}: \mathcal{D}^* \rightarrow \mathcal{R}$ be a randomized algorithm that satisfies (ϵ, δ) -DP with $\epsilon \in \mathbb{R}_+$ and $\delta \in [0, 1]$ if for all neighboring datasets $D \sim D'$, and for all possible subsets $R \subseteq \mathcal{R}$ of the result space \mathcal{M} satisfies

$$\mathbb{P}[\mathcal{M}(D) \in R] \leq e^\epsilon \cdot \mathbb{P}[\mathcal{M}(D') \in R] + \delta$$

Fairness Objective: Demographic Parity

Definition (Demographic Disparity)

$$\Gamma_{\text{DemParity}}(k, z) = \mathbb{P}[\hat{Y} = k' \mid Z = z] - \mathbb{P}[\hat{Y} = k' \mid Z \neq z]$$

where $\hat{Y} = \omega(\mathbf{x}, z)$ are model $\omega : \mathcal{X} \times \mathcal{Z} \mapsto \mathcal{K}$ predictions for samples with sensitive attribute z .

Definition (γ -disparity)

$\forall z \in \mathcal{Z}, \forall k \in \mathcal{K},$

$$\Gamma_{\text{DemParity}}(k, z) \leq \gamma$$

ML trustworthiness as multi-objective optimization

$$\begin{aligned} & \min_{\omega} && l_{\text{acc}}(\omega) \\ & \text{subject to} && l_{\text{priv}}(\omega) \leq \varepsilon \\ & && l_{\text{fair}}(\omega) \leq \gamma \end{aligned} \tag{1}$$

where $(l_{\text{acc}}, l_{\text{priv}}, l_{\text{fair}}) \in \mathbb{R}_{\geq 0}^3$ are the loss functions for each of the utility, privacy and fairness criteria, respectively.

ML trustworthiness as multi-objective optimization

$$\begin{aligned} & \min_{\omega} && l_{\text{acc}}(\omega) \\ & \text{subject to} && l_{\text{priv}}(\omega) \leq \varepsilon \\ & && l_{\text{fair}}(\omega) \leq \gamma \end{aligned} \tag{1}$$

where $(l_{\text{acc}}, l_{\text{priv}}, l_{\text{fair}}) \in \mathbb{R}_{\geq 0}^3$ are the loss functions for each of the utility, privacy and fairness criteria, respectively.

- Problem 1: Privacy is ensured at the level of mechanism (here, the ML pipeline)

ML trustworthiness as multi-objective optimization

$$\begin{aligned} & \min_{\omega} && l_{\text{acc}}(\omega) \\ & \text{subject to} && \cancel{l_{\text{priv}}(\omega) \leq \epsilon} \\ & && l_{\text{fair}}(\omega) \leq \gamma \end{aligned} \tag{1}$$

where $(l_{\text{acc}}, l_{\text{priv}}, l_{\text{fair}}) \in \mathbb{R}_{\geq 0}^3$ are the loss functions for each of the utility, privacy and fairness criteria, respectively.

- Problem 1: Privacy is ensured at the level of mechanism (here, the ML pipeline) \Rightarrow We do not have a sample-based privacy loss

ML trustworthiness as multi-objective optimization

Treated as hyper-parameters

$$\begin{aligned} & \min_{\omega} && l_{\text{acc}}(\omega) & | \\ \text{subject to} &&& \cancel{l_{\text{priv}}(\omega) \leq \epsilon} & \\ &&& l_{\text{fair}}(\omega) \leq \gamma & \end{aligned} \quad (1)$$

where $(l_{\text{acc}}, l_{\text{priv}}, l_{\text{fair}}) \in \mathbb{R}_{\geq 0}^3$ are the loss functions for each of the utility, privacy and fairness criteria, respectively.

- Problem 1: Privacy is ensured at the level of mechanism (here, the ML pipeline) \Rightarrow We do not have a sample-based privacy loss
- Problem 2: Trustworthy parameters are treated as hyper-parameters, not first-class objectives

ML trustworthiness as multi-objective optimization

Treated as hyper-parameters

$$\begin{array}{ll} \min_{\omega} & l_{\text{acc}}(\omega) \\ \text{subject to} & l_{\text{priv}}(\omega) \leq \epsilon \\ & l_{\text{fair}}(\omega) \leq \gamma \end{array} \quad (1)$$

where $(l_{\text{acc}}, l_{\text{priv}}, l_{\text{fair}}) \in \mathbb{R}_{\geq 0}^3$ are the loss functions for each of the utility, privacy and fairness criteria, respectively.

- Problem 1: Privacy is ensured at the level of mechanism (here, the ML pipeline) \Rightarrow We do not have a sample-based privacy loss
- Problem 2: Trustworthy parameters are treated as hyper-parameters, not first-class objectives \Rightarrow Pre-Selection Bias

Pre-selection Bias

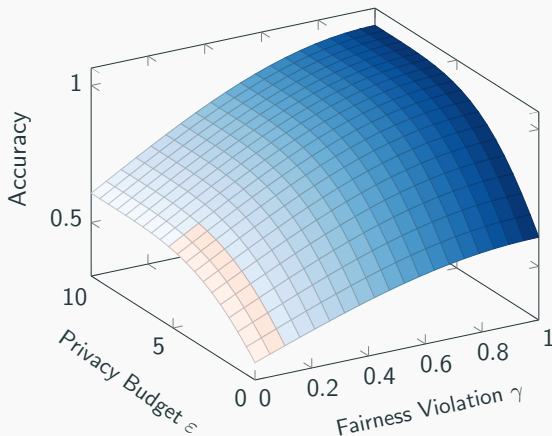
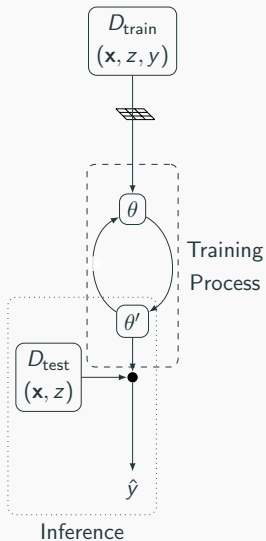
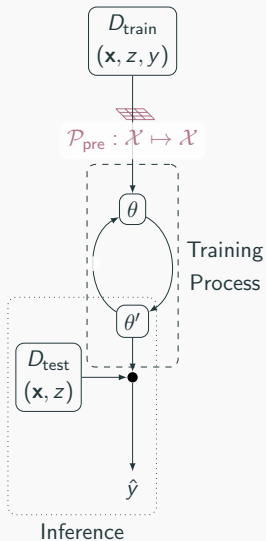


Figure 2: Pre-selection of trustworthiness parameters only recovers a portion of the Pareto frontier. The remaining parts of frontier (shaded blue) are never explored.

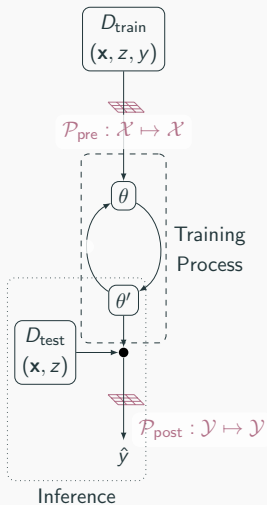
DP Learning to Fair DP Learning: Fairness Intervention



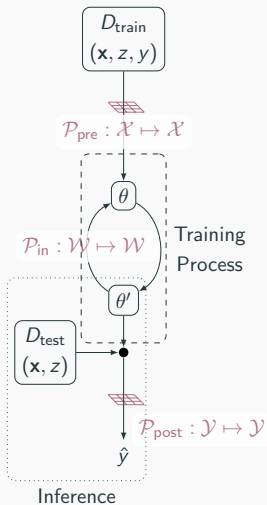
DP Learning to Fair DP Learning: Fairness Intervention



DP Learning to Fair DP Learning: Fairness Intervention



DP Learning to Fair DP Learning: Fairness Intervention

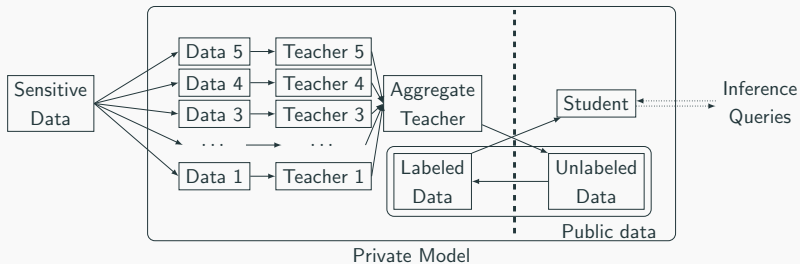


Theorem

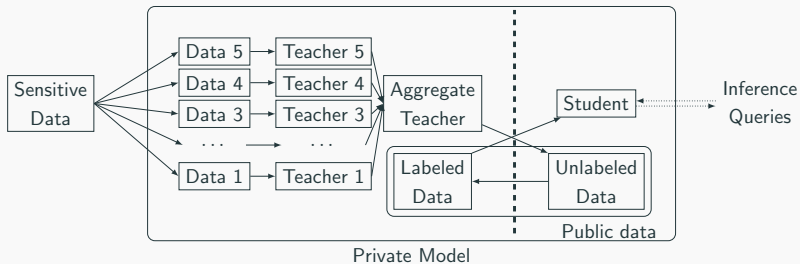
Assume the training dataset

$D = \{(\mathbf{x}, z, y) \mid \mathbf{x} \in \mathcal{X}, z \in \mathcal{Z}, y \in \mathcal{Y}\}$ is fed through the demographic parity pre-processor \mathcal{P}_{pre} following an ordering defined over the input space \mathcal{X} . Let \mathcal{P}_{pre} enforce a maximum violation γ , and $|\mathcal{Z}| = 2$. Suppose now \mathcal{M} is an (ϵ, δ) training mechanism, then $\mathcal{M} \circ \mathcal{P}_{pre}$ is $(K_\gamma \epsilon, K_\gamma e^{K_\gamma \epsilon} \delta)$ -DP where $K_\gamma = 2 + \left\lceil \frac{2\gamma}{1-\gamma} \right\rceil$.

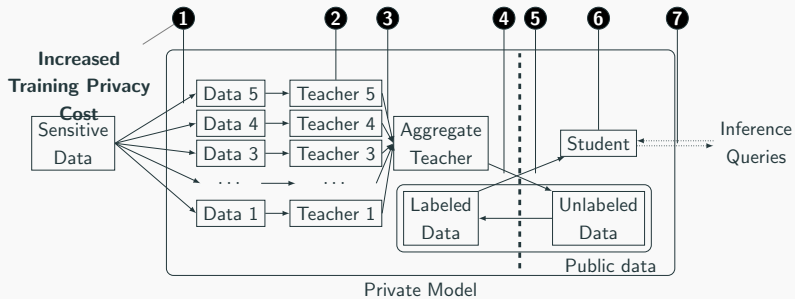
Example: PATE



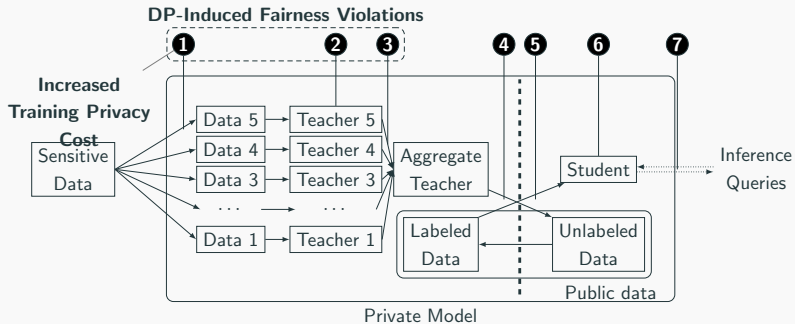
Example: PATE



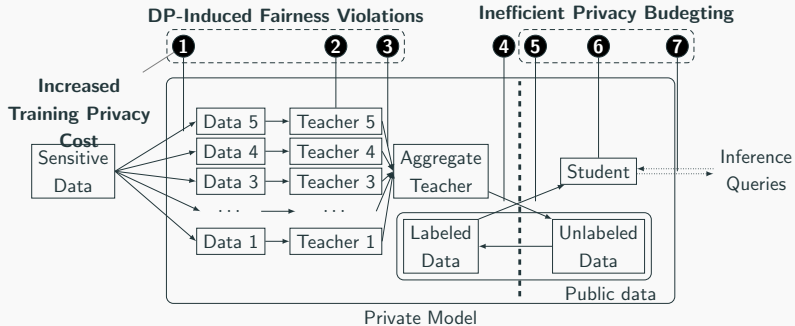
Example: PATE



Example: PATE



Example: PATE



Algorithm 1 Confident-GNMax Aggregator

Input: query data point x , sensitive attribute z , predicted class label k , subpopulation subclass counts $m : \mathcal{Z} \times \mathcal{K} \mapsto \mathbb{Z}_{\geq 0}$

Require: minimum count M , threshold T , noise parameters σ_1, σ_2 , fairness violation margin γ

- 1: **if** $\max_j \{n_j(x)\} + \mathcal{N}(0, \sigma_1^2) \geq T$ **then**
- 2: $k \leftarrow \arg \max_j \{n_j(x) + \mathcal{N}(0, \sigma_2^2)\}$
- 3: **return** k

13: **else**

14: **return** \perp

Algorithm 2 Confident&Fair-GNMax Aggregator

Input: query data point x , sensitive attribute z , predicted class label k , subpopulation subclass counts $m : \mathcal{Z} \times \mathcal{K} \mapsto \mathbb{Z}_{\geq 0}$

Require: minimum count M , threshold T , noise parameters σ_1, σ_2 , fairness violation margin γ

1: **if** $\max_j \{n_j(x)\} + \mathcal{N}(0, \sigma_1^2) \geq T$ **then**

2: $k \leftarrow \arg \max_j \{n_j(x) + \mathcal{N}(0, \sigma_2^2)\}$

3: **if** $\sum_{\tilde{k}} m(z, \tilde{k}) < M$ **then**

4: $m(z, k) \leftarrow m(z, k) + 1$

5: **return** k

12: **else**

13: **return** \perp

Algorithm 2 Confident&Fair-GNMax Aggregator

Input: query data point x , sensitive attribute z , predicted class label k , subpopulation subclass counts $m : \mathcal{Z} \times \mathcal{K} \mapsto \mathbb{Z}_{\geq 0}$

Require: minimum count M , threshold T , noise parameters σ_1, σ_2 , fairness violation margin γ

```
1: if  $\max_j \{n_j(x)\} + \mathcal{N}(0, \sigma_1^2) \geq T$  then
2:    $k \leftarrow \arg \max_j \{n_j(x) + \mathcal{N}(0, \sigma_2^2)\}$ 
3:   if  $\sum_{\tilde{k}} m(z, \tilde{k}) < M$  then
4:      $m(z, k) \leftarrow m(z, k) + 1$ 
5:     return  $k$ 
6:   else
7:     if  $\left( \frac{m(z, k) + 1}{(\sum_{\tilde{k}} m(z, \tilde{k})) + 1} - \frac{\sum_{\tilde{z} \neq z} m(\tilde{z}, k)}{\sum_{\tilde{z} \neq z, \tilde{k}} m(\tilde{z}, \tilde{k})} \right) < \gamma$  then
8:        $m(z, k) \leftarrow m(z, k) + 1$ 
9:       return  $k$ 
10:    else
11:      return  $\perp$ 
12:  else
13:    return  $\perp$ 
```

Closing the fairness gap with Reject-Option for Fairness

- Optimizing for fairness during the training process does not guarantee that fairness is obtained at inference time
- What if there were a **hard constraint** on fairness violations at inference time?
- A **reject-option** allows to refuse to answer a query at inference time for fairness purposes.
- Introduces a new utility dimension:
Coverage := $\frac{\# \text{ Queries Answered}}{\# \text{ Queries}}$

Algorithm 5 Inference-time Demographic Parity Post-Processor (IDP³)

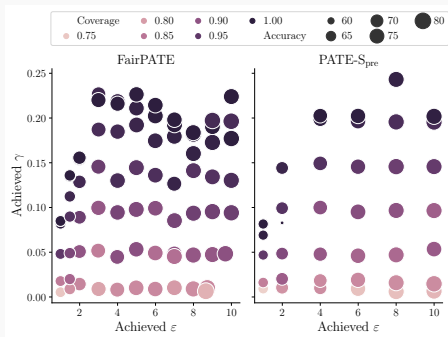
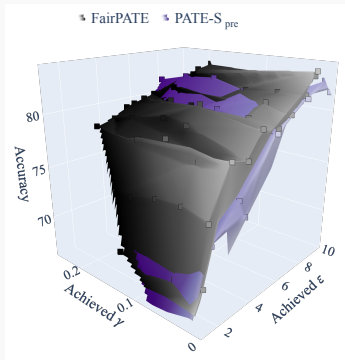
Input: data point x , sensitive attribute z , predicted label \hat{y} ,
subpopulation-class counts $m : \mathcal{Z} \times \mathcal{Y} \mapsto \mathbb{Z}_{\geq 0}$

Require: minimum count M , fairness violation margin γ

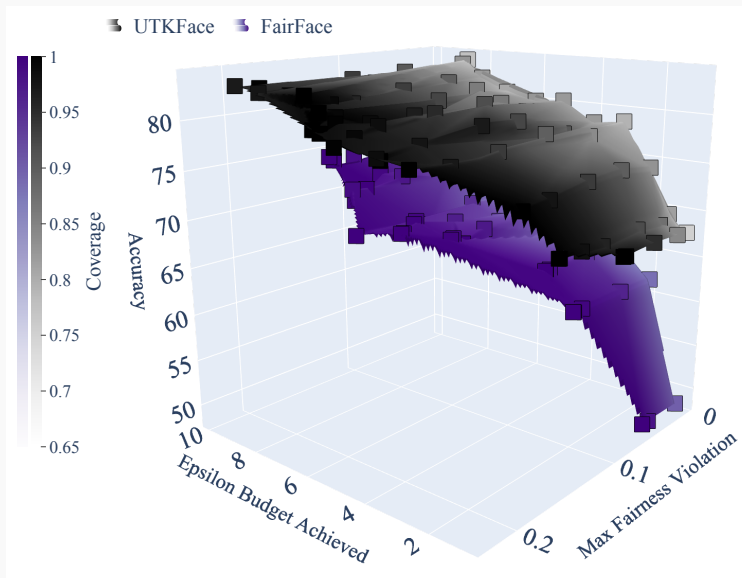
- 1: **if** $\sum_{\tilde{y}} m(z, \tilde{y}) < M$ **then**
 - 2: $m(z, y) \leftarrow m(z, \hat{y}) + 1$
 - 3: **return** \hat{y}
 - 4: **else**
 - 5: **if** $\left(\frac{m(z, \hat{y}) + 1}{(\sum_{\tilde{y}} m(z, \tilde{y})) + 1} - \frac{\sum_{\tilde{z} \neq z} m(\tilde{z}, \hat{y})}{\sum_{\tilde{z} \neq z, \tilde{y}} m(\tilde{z}, \tilde{y})} \right) < \gamma$ **then**
 - 6: $m(z, y) \leftarrow m(z, \hat{y}) + 1$
 - 7: **return** \hat{y}
 - 8: **else**
 - 9: **return** \perp
-

Results

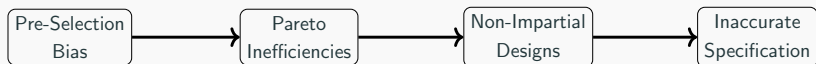
FairPATE Pareto-dominates similar designs in most contexts



Specification without direct data access is possible



Conclusion: Specification requires objective-impartiality!



Conclusion: Specification requires objective-impartiality!



Conclusion: Specification requires objective-impartiality!



Conclusion: Specification requires objective-impartiality!



Conclusion: Specification requires objective-impartiality!



Thank you!