

Impacts of Data Privacy and Equity on Public Policy

Ferdinando Fioretto
Syracuse U.

Claire McKay Bowen
Urban Institute

FAccT 2022

Part I

*What is data privacy and confidentiality and
why is it important?*



BRIEFING ROOM

Executive Order On Advancing Racial Equity and Support for Underserved Communities Through the Federal Government

JANUARY 20, 2021 • PRESIDENTIAL ACTIONS

“...advanc[e] equity for all, including people of color and others who have been historically underserved, marginalized, and adversely affected by persistent poverty and inequality.”



ONE NATION, TRACKED

AN INVESTIGATION INTO THE SMARTPHONE TRACKING
INDUSTRY FROM TIMES OPINION

Lack of Public Data Hampers COVID-19 Fight

STATELINE ARTICLE August 3, 2020 By: Christine Vestal Topics: Health Read time: 7 min



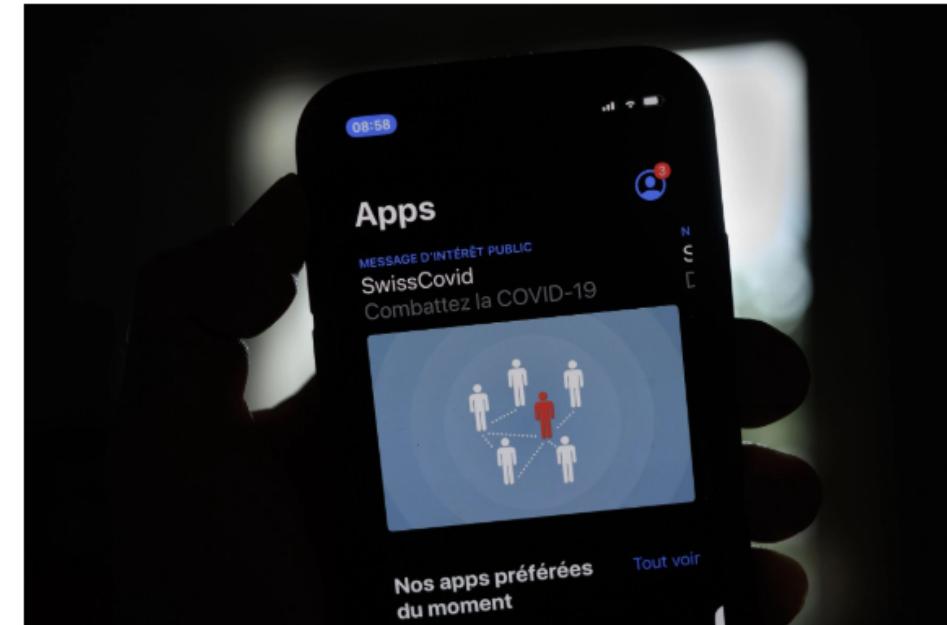
A ventilator helps a COVID-19 patient breathe at a Houston hospital. Hospital data related to the coronavirus pandemic will now be collected by a private technology firm, rather than the Centers for Disease Control and Prevention. Epidemiologists say better COVID-19 data is needed to improve the nation's response.

David J. Phillip/The Associated Press

• U R B A N • I N S T I T U T E •

Google Promises Privacy With Virus App but Can Still Collect Location Data

Some government agencies that use the software said they were surprised that Google may pick up the locations of certain app users. Others said they had unsuccessfully pushed Google to make a change.



Switzerland has asked Google to decouple the location setting requirement on Android phones from Bluetooth, which the country's virus alert app uses to detect nearby smartphones. Fabrice Coffrini/Agence France-Presse — Getty Images

By Natasha Singer

July 20, 2020



'It's not a pretty picture': Why the lack of racial data around COVID vaccines is 'massive barrier' to better distribution

Nada Hassanein USA TODAY

Published 5:30 a.m. ET Feb. 1, 2021 | Updated 2:10 p.m. ET Feb. 1, 2021



Abigail Echo-Hawk, chief research officer with Seattle Indian Health Board and a member of the Pawnee Tribe, gets a shot of the Moderna COVID-19 vaccine on Dec. 21. A colleague used a black pen to inscribe "For the (Heart) love of Native People" over the injection spot. Karen Ducey, Getty Images

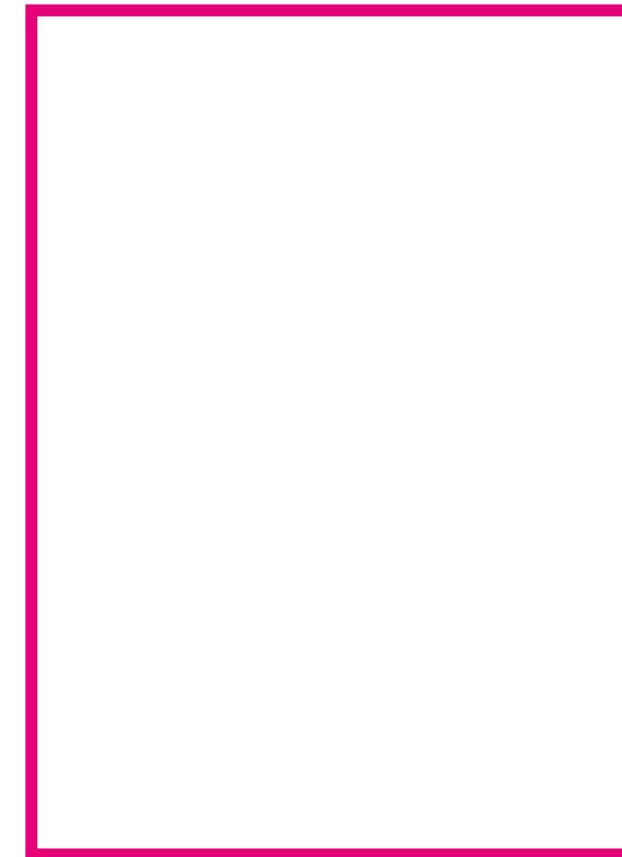
*How do public policymakers gain access to
confidential data?*

How do we access confidential data?

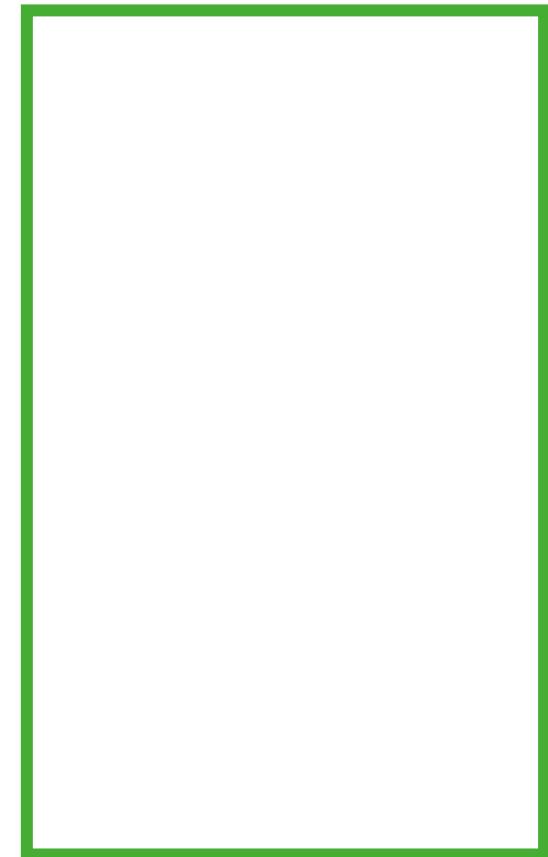
User Interface
Layer



Privacy
Layer

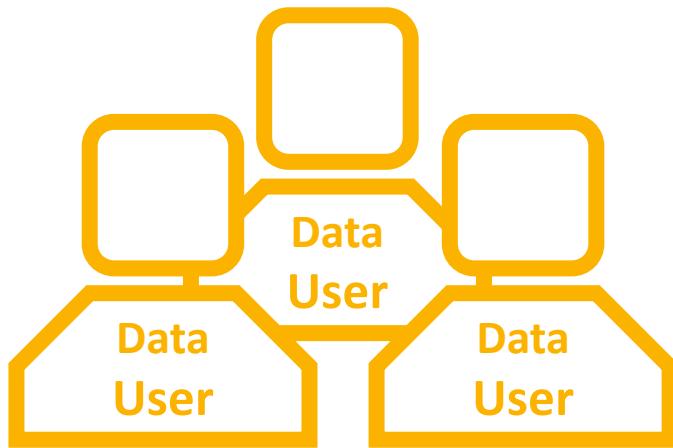


Data Access Layer

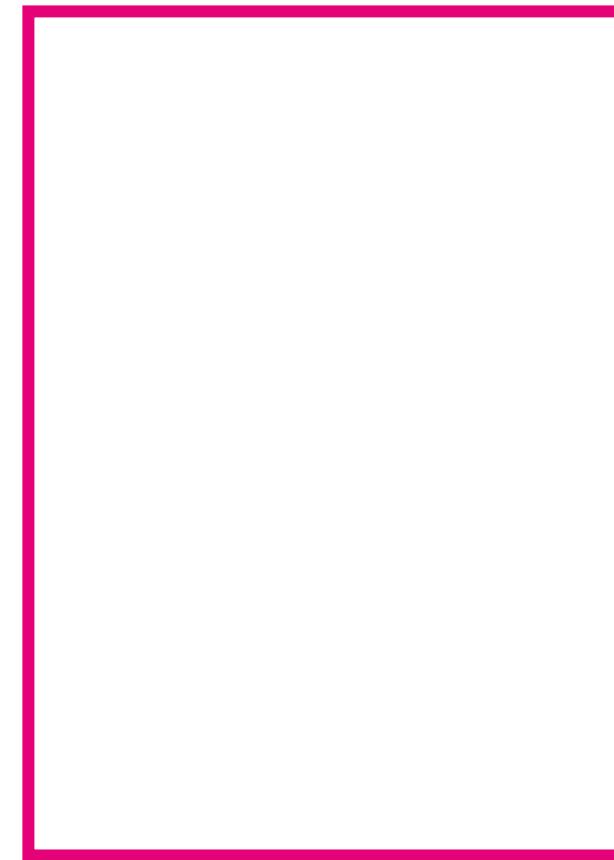


How do we access confidential data?

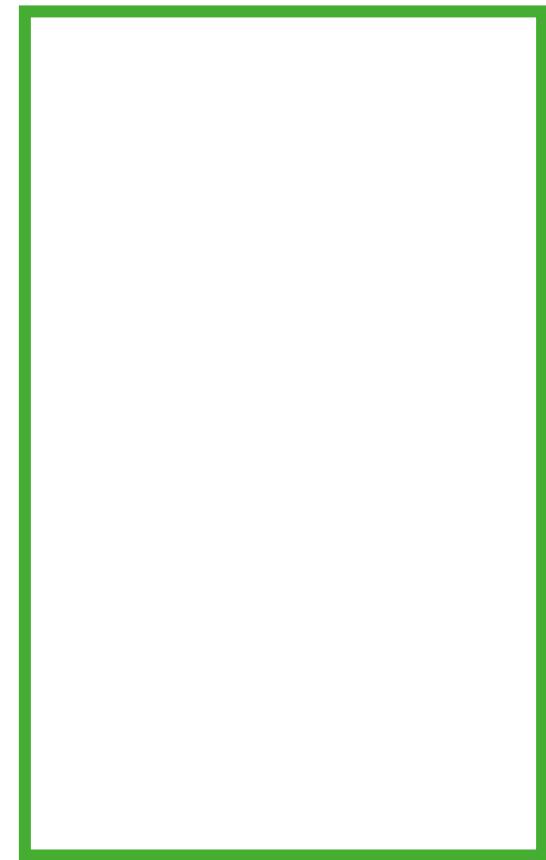
User Interface
Layer



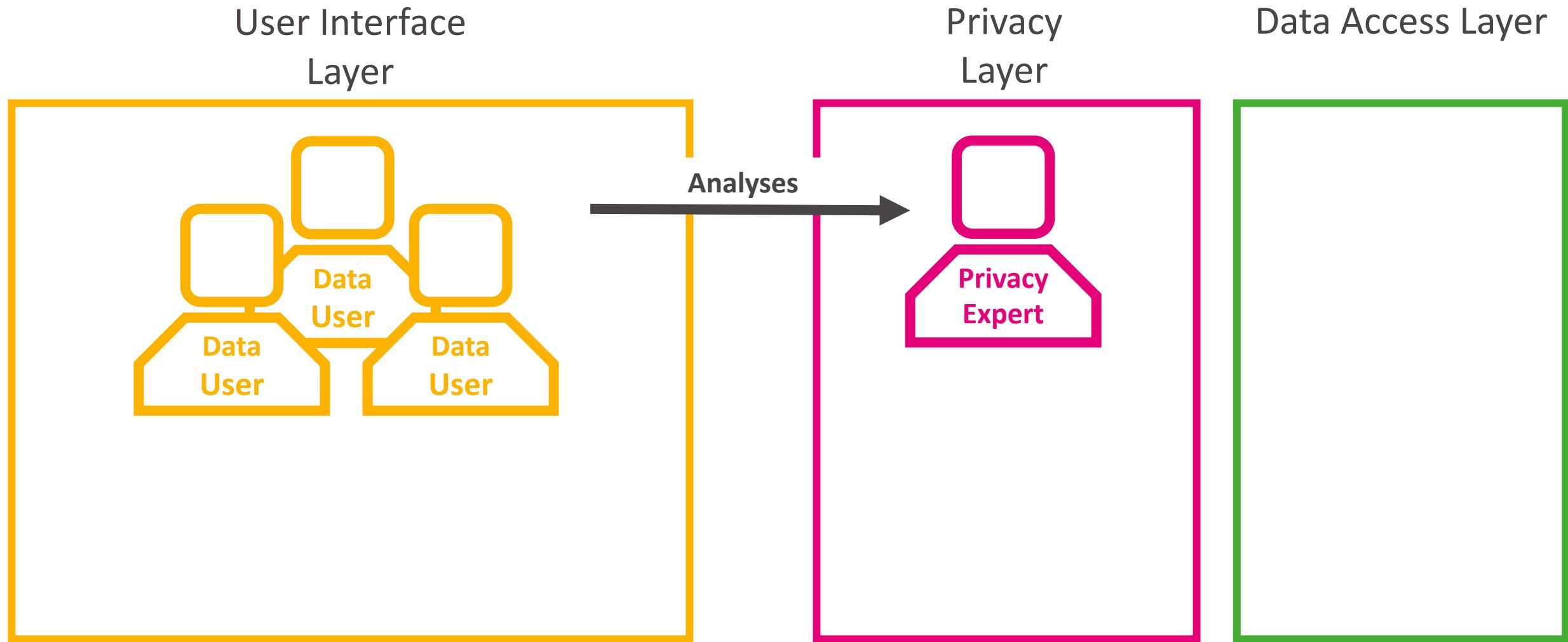
Privacy
Layer



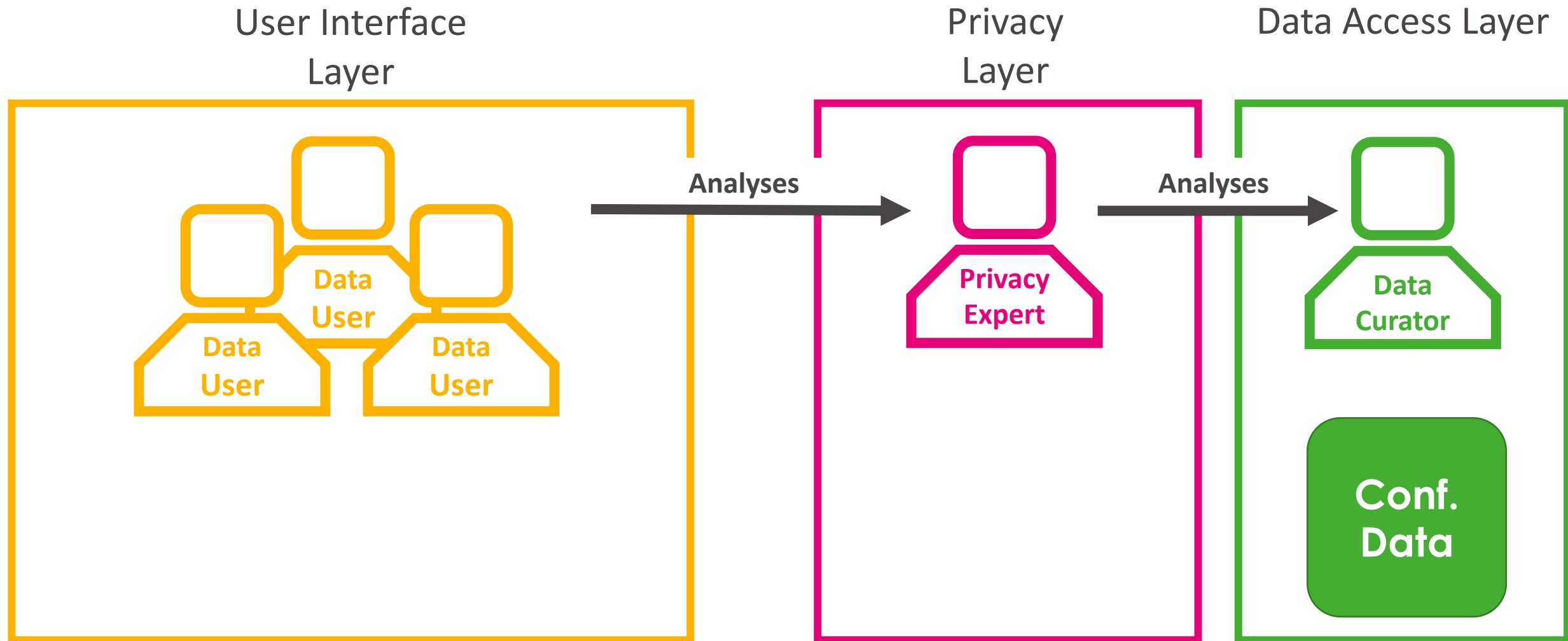
Data Access Layer



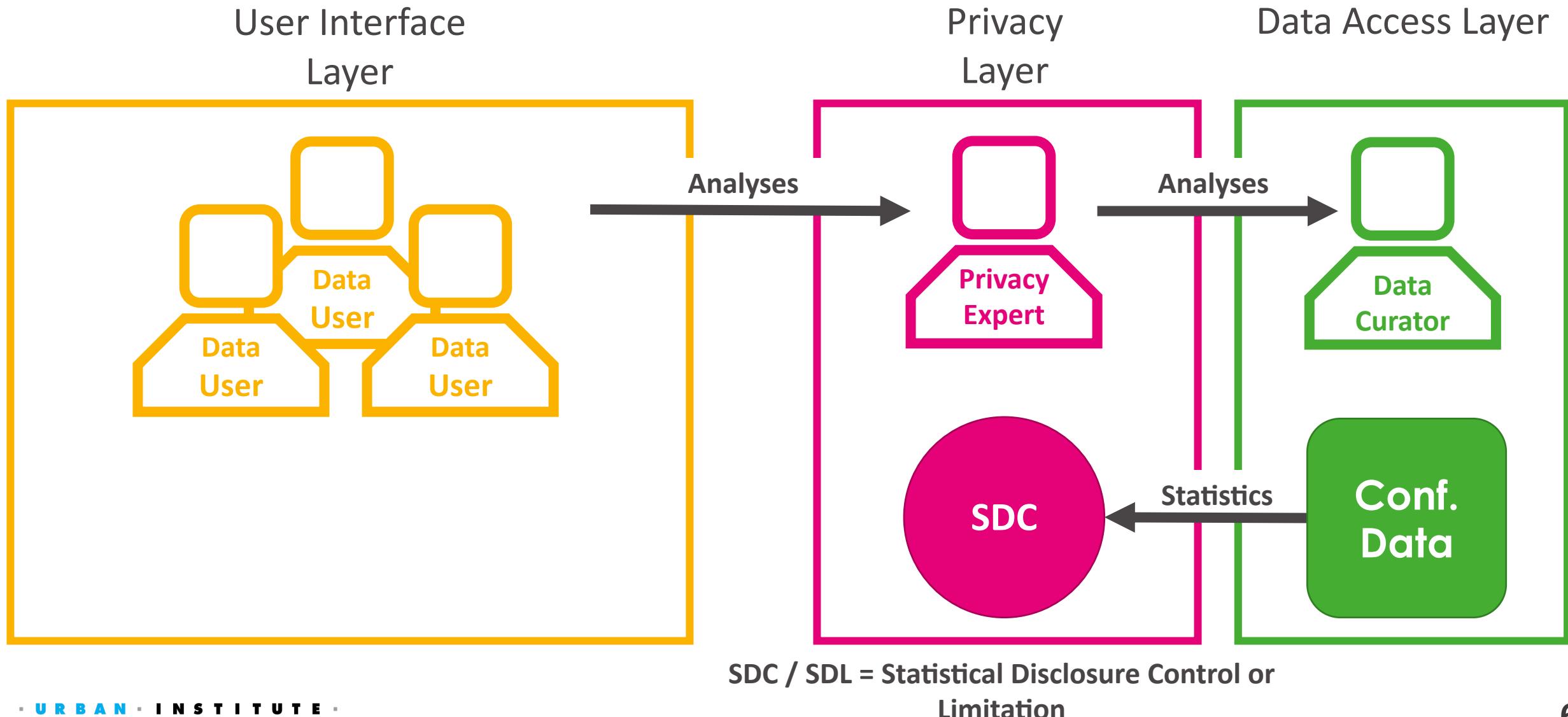
How do we access confidential data?



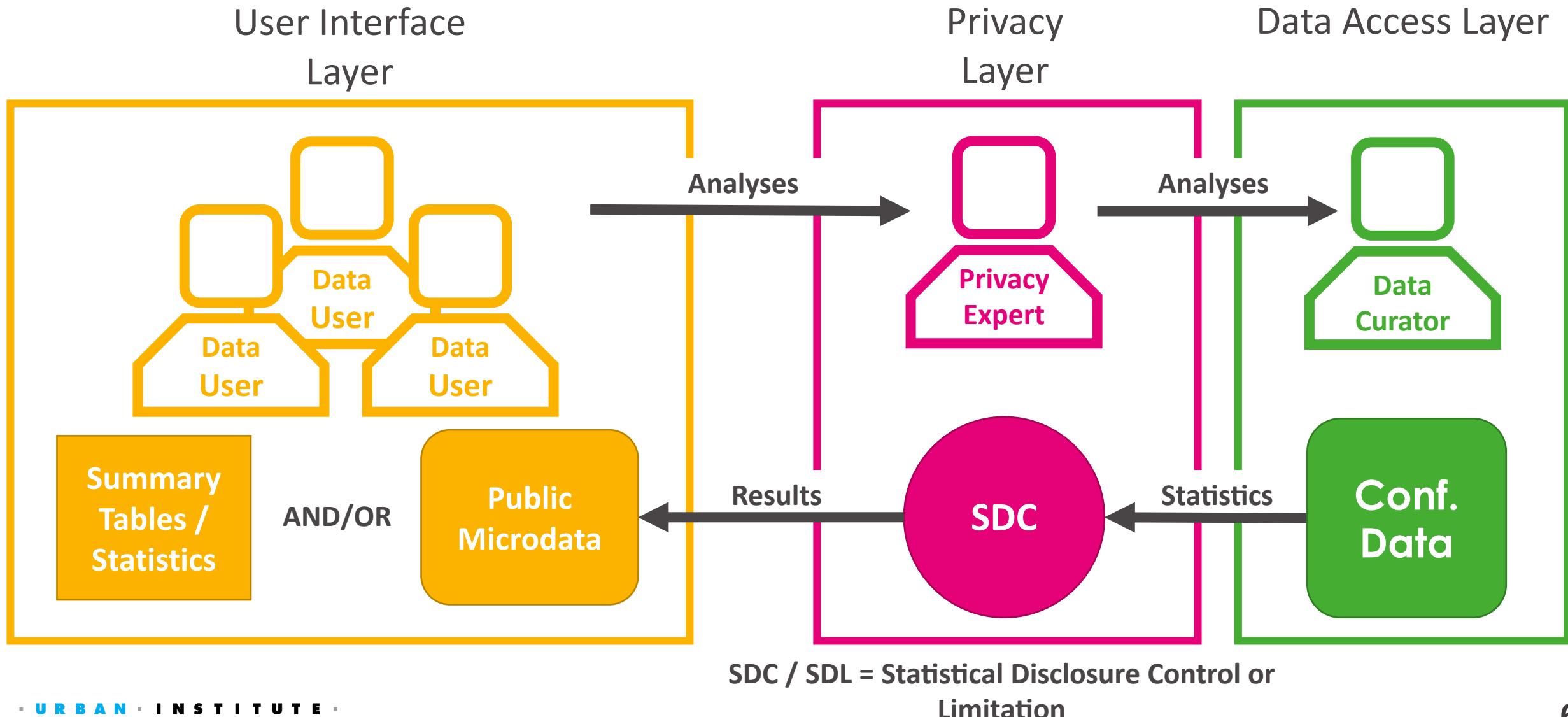
How do we access confidential data?



How do we access confidential data?

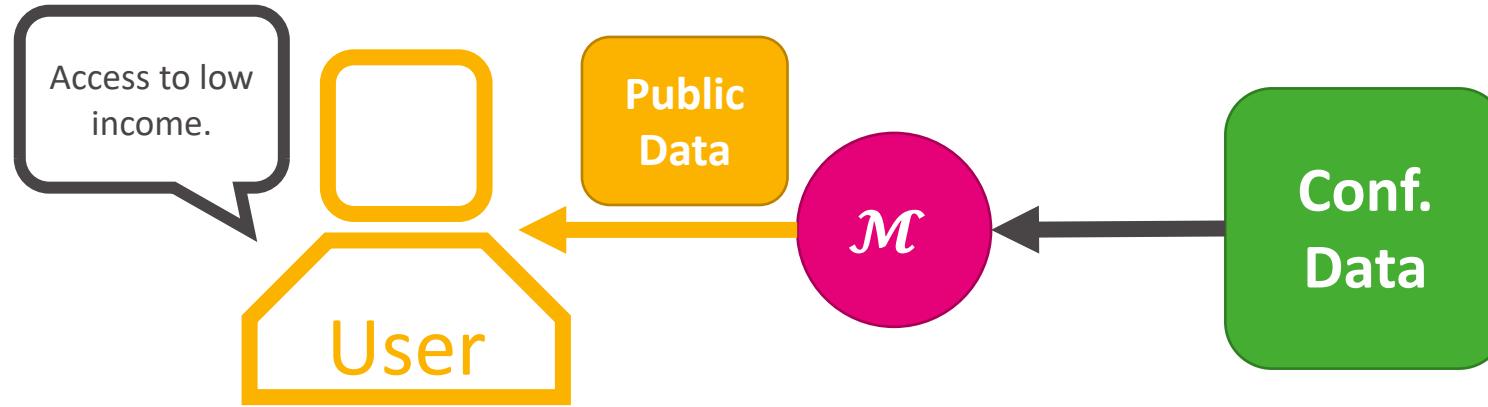


How do we access confidential data?

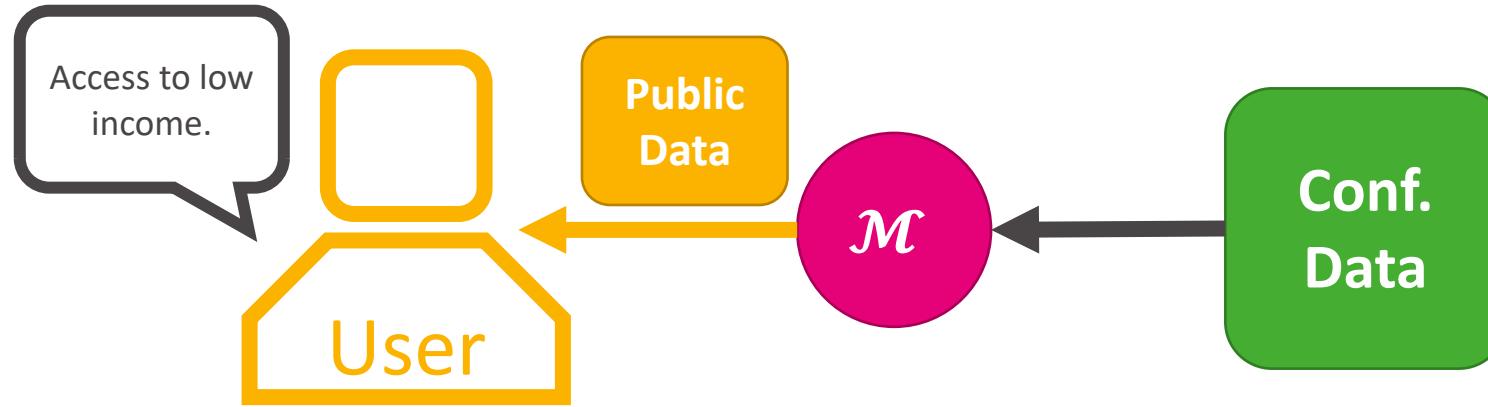


Why is it hard to implementing this framework?

Why is it hard to implementing this framework?

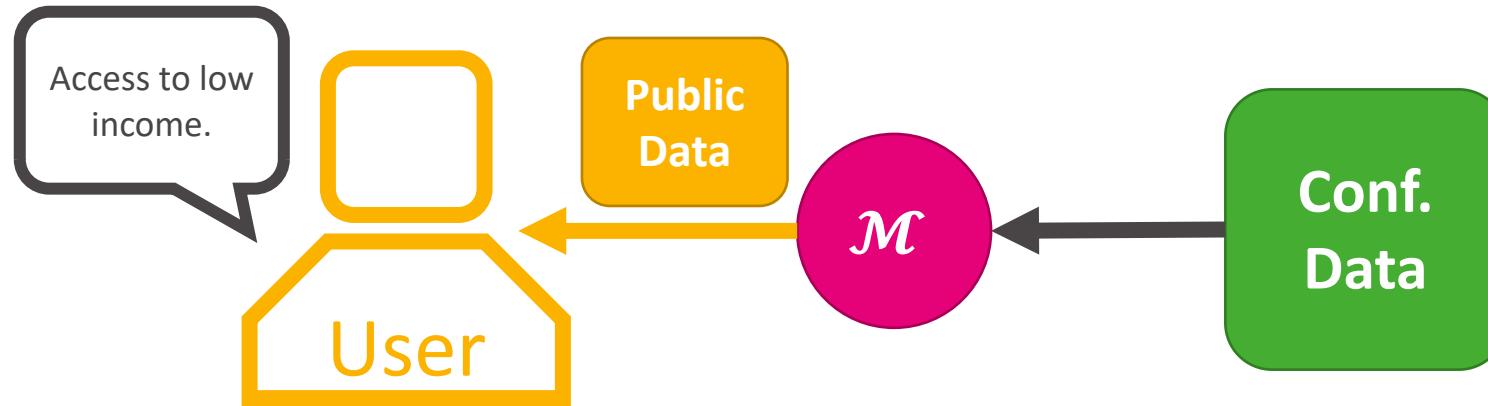


Why is it hard to implementing this framework?

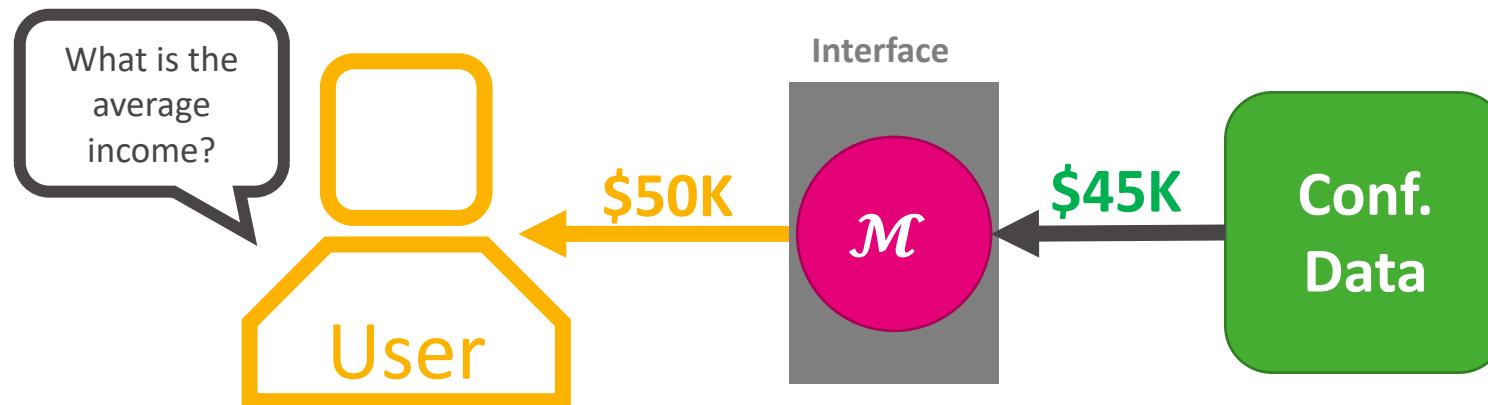


How do we **measure utility** (usefulness) and **disclosure risk** of the data?

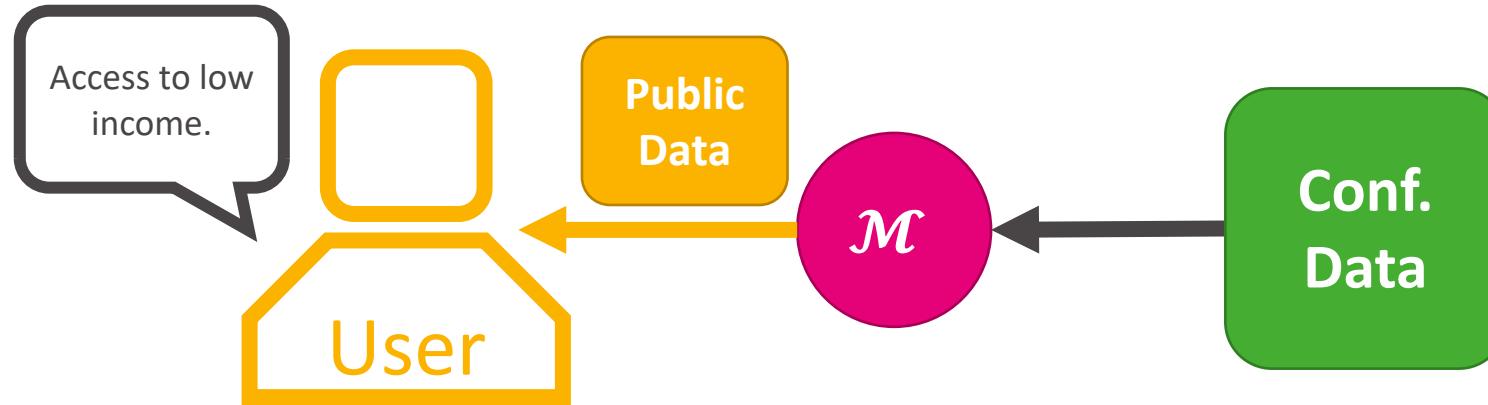
Why is it hard to implementing this framework?



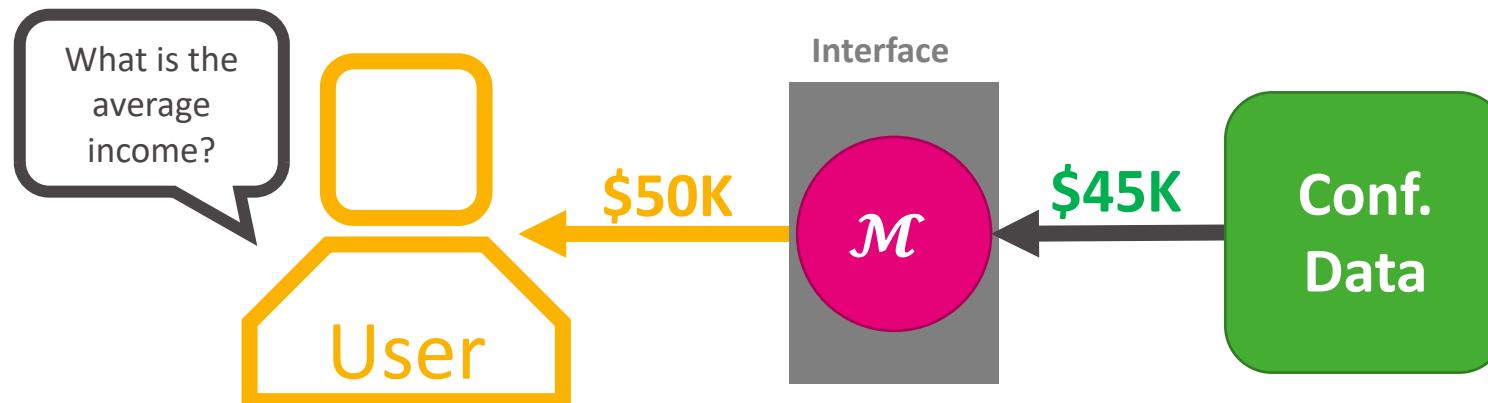
How do we **measure utility** (usefulness) and **disclosure risk** of the data?



Why is it hard to implementing this framework?



How do we **measure utility** (usefulness) and **disclosure risk** of the data?

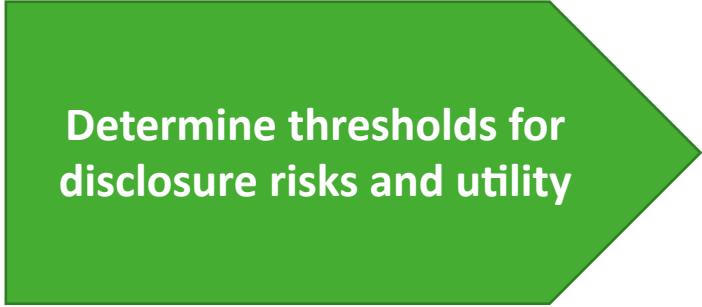


How **much noise should be added** and how do you **limit the number of queries**?

What is the workflow?

What is the workflow?

Step 1



Determine thresholds for disclosure risks and utility

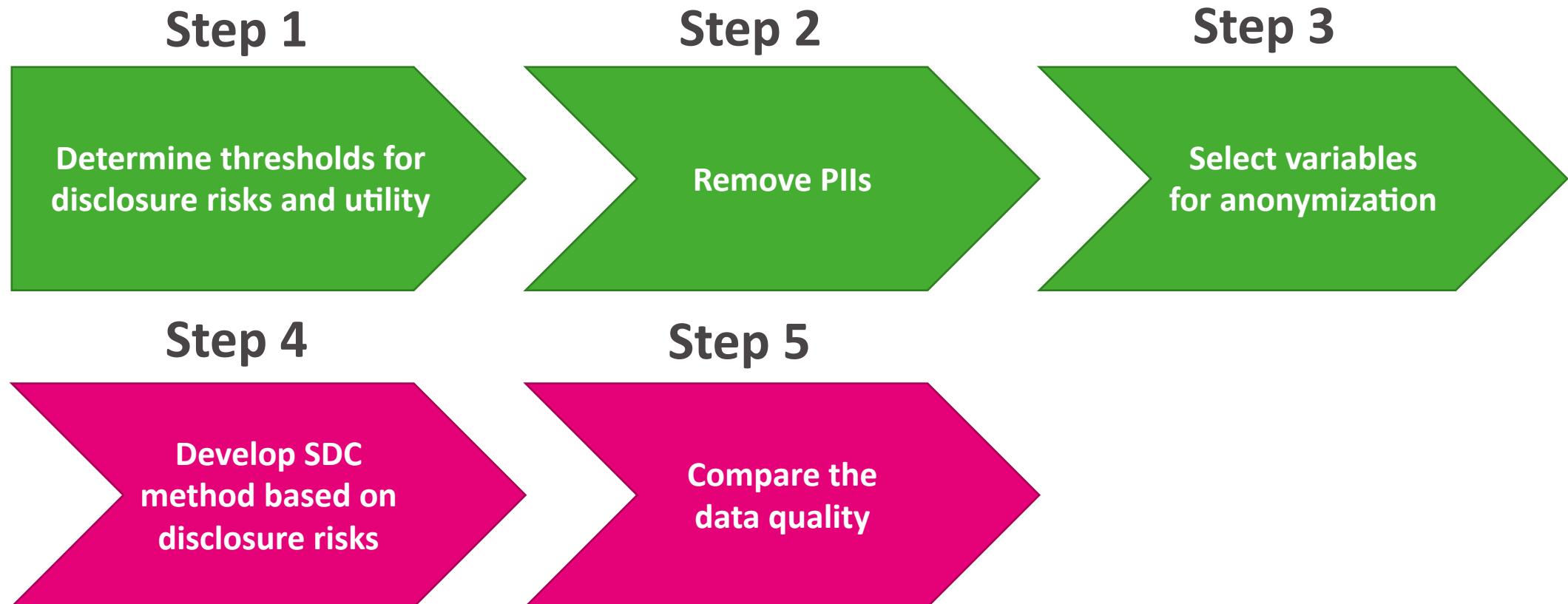
What is the workflow?



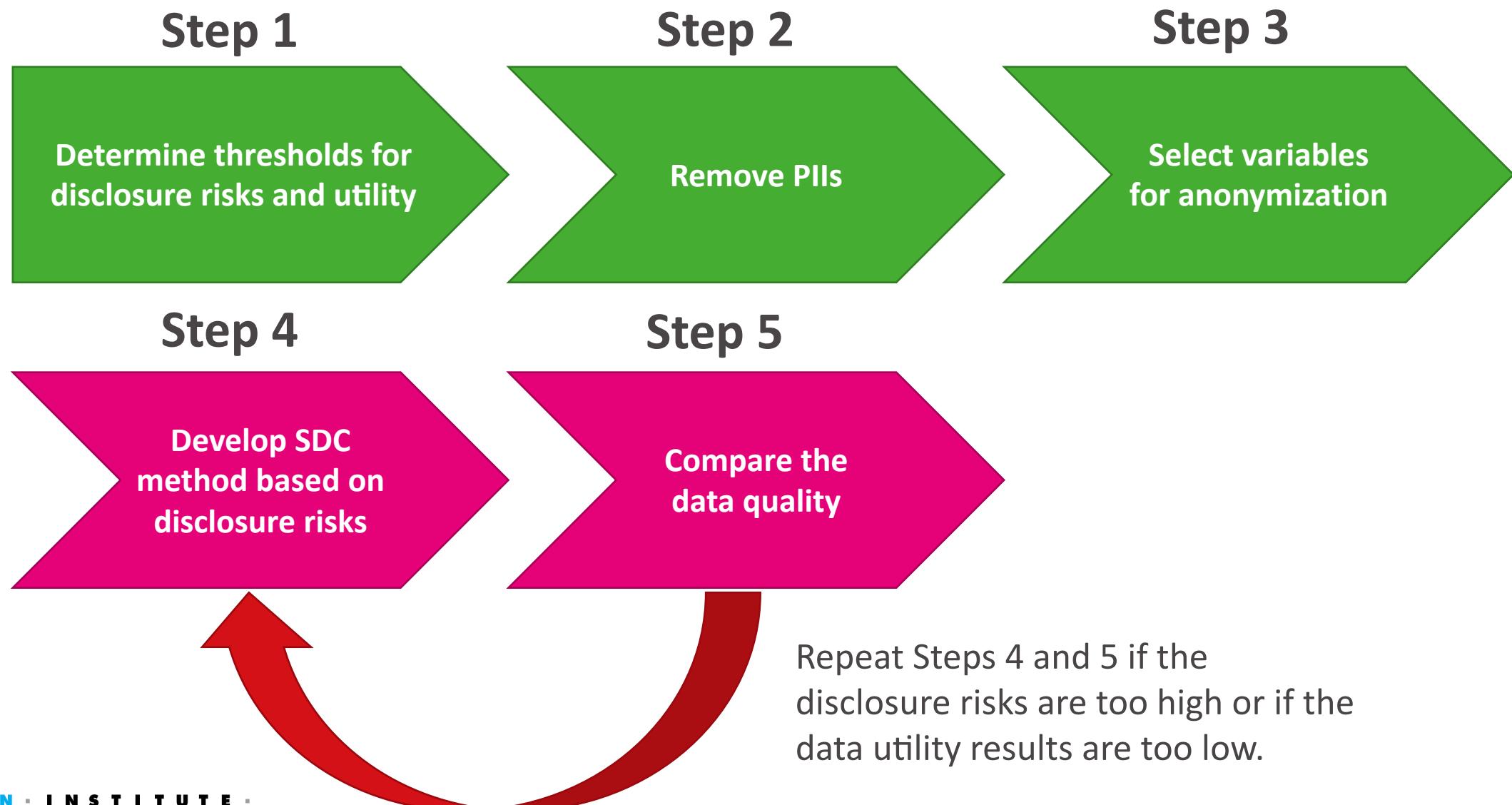
What is the workflow?



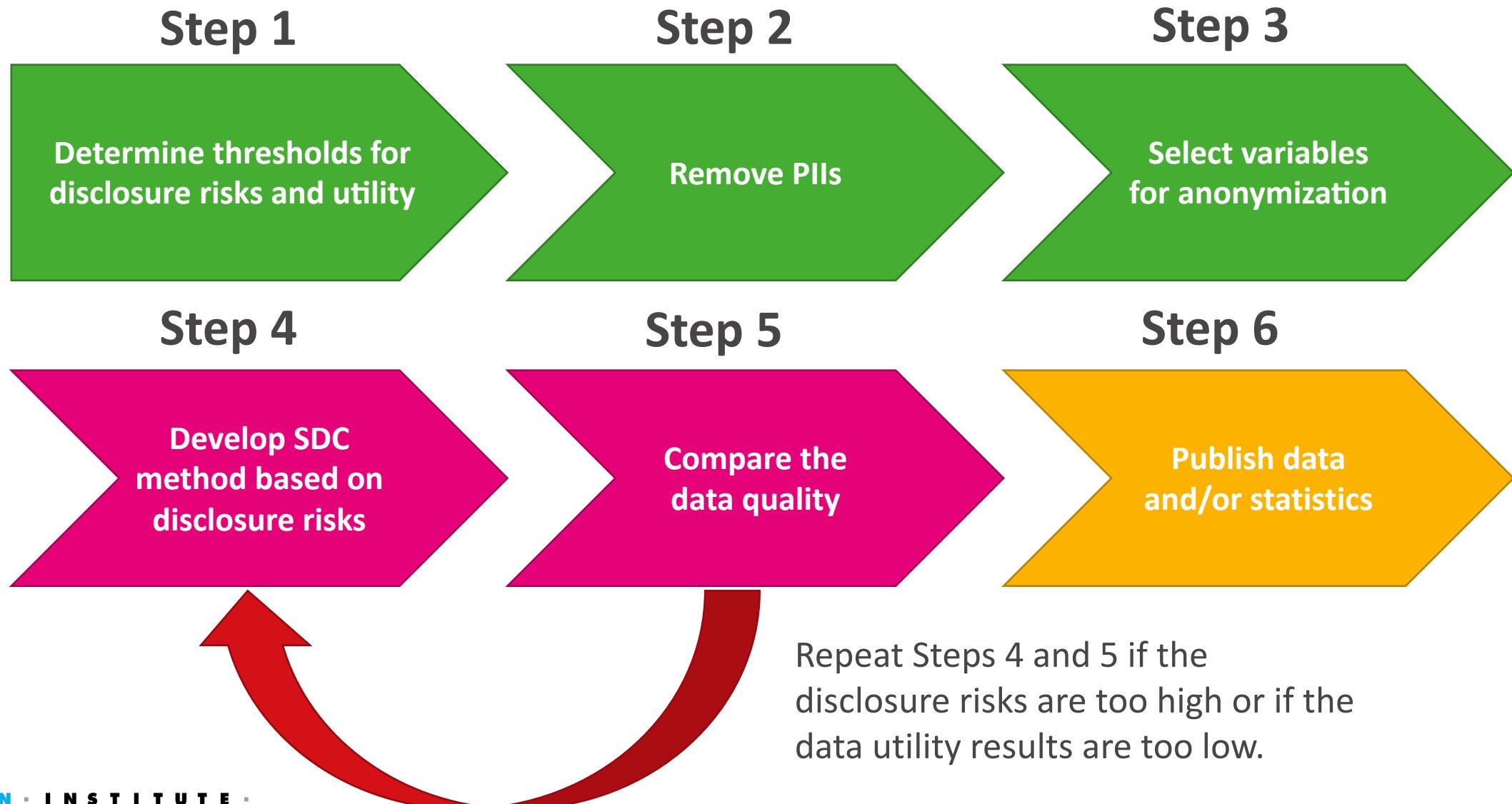
What is the workflow?



What is the workflow?



What is the workflow?



How are these methods structured?

How are these methods structured?

- 1. Pre-Processing Step:** determining priorities for which statistics to preserve

How are these methods structured?

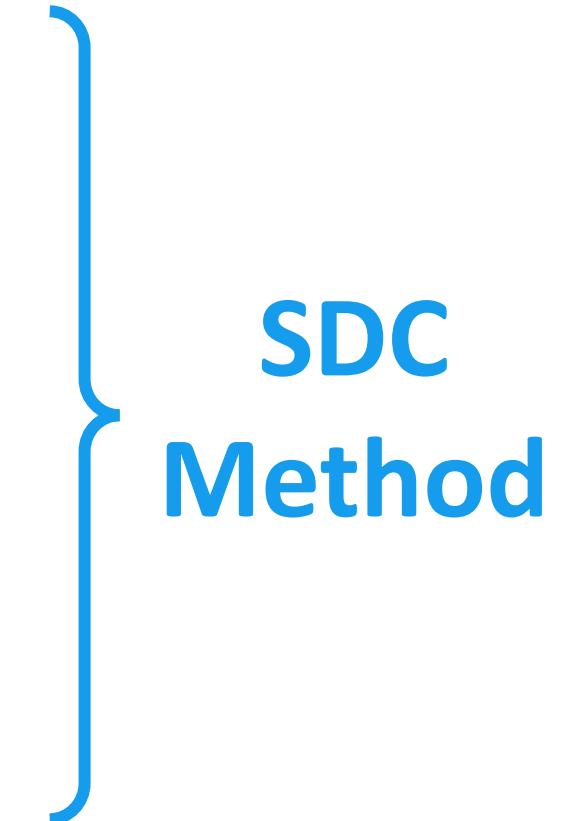
- 1. Pre-Processing Step:** determining priorities for which statistics to preserve
- 2. Privacy Step:** applying a sanitizer to the desired statistic

How are these methods structured?

- 1. Pre-Processing Step:** determining priorities for which statistics to preserve
- 2. Privacy Step:** applying a sanitizer to the desired statistic
- 3. Post-Processing Step:** ensuring the results of the statistics are consistent with realistic constraints, such as negative population counts

How are these methods structured?

- 1. Pre-Processing Step:** determining priorities for which statistics to preserve
- 2. Privacy Step:** applying a sanitizer to the desired statistic
- 3. Post-Processing Step:** ensuring the results of the statistics are consistent with realistic constraints, such as negative population counts



How are data used for public policy decision making?

Economic Mobility – U.S. Taxpayer Data



School District Boundaries – U.S. Education Data

Dividing Lines

How School Districts Draw Attendance Boundaries to
Perpetuate School Segregation

September 14, 2021

2020 Census Data

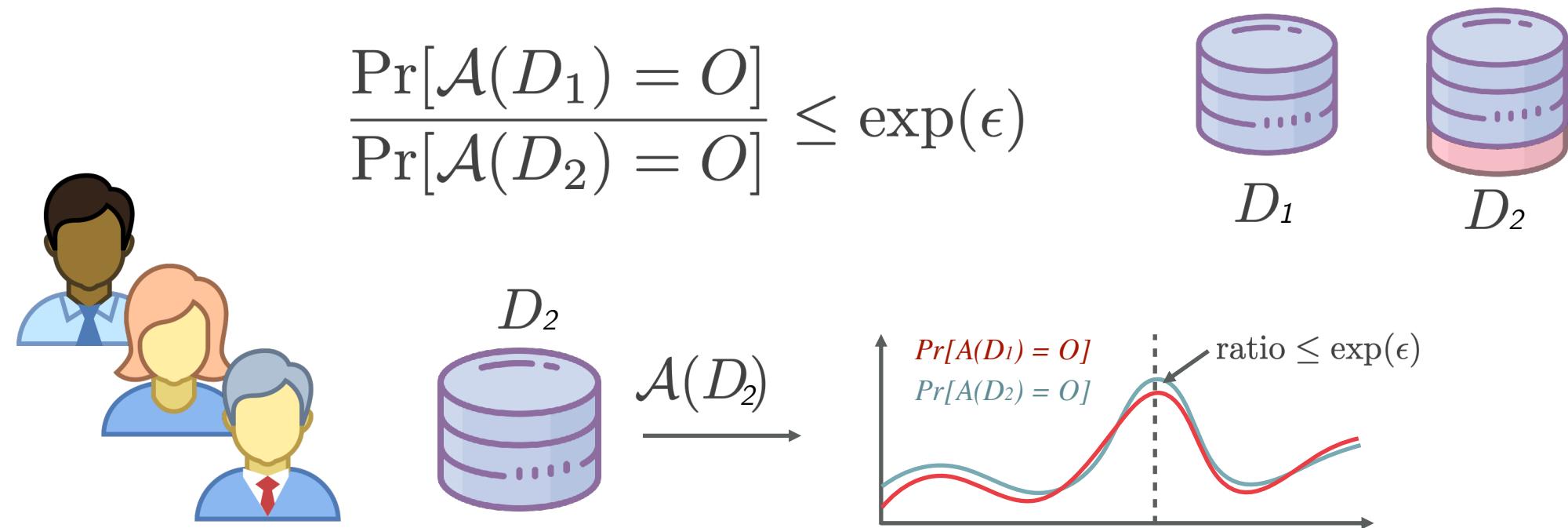
- Allocating \$1.5 trillion budget
- Apportioning 435 congressional seats
- Planning for natural disasters
- Understanding economic well-being
- Determining the number of restaurant permits



Part II

Differential Privacy

A randomized algorithm \mathcal{A} is ϵ -differentially private if, for all pairs of inputs D_1, D_2 , differing in one entry, and for any output O :



Intuition: An adversary should not be able to use output O to distinguish between any D_1 and D_2

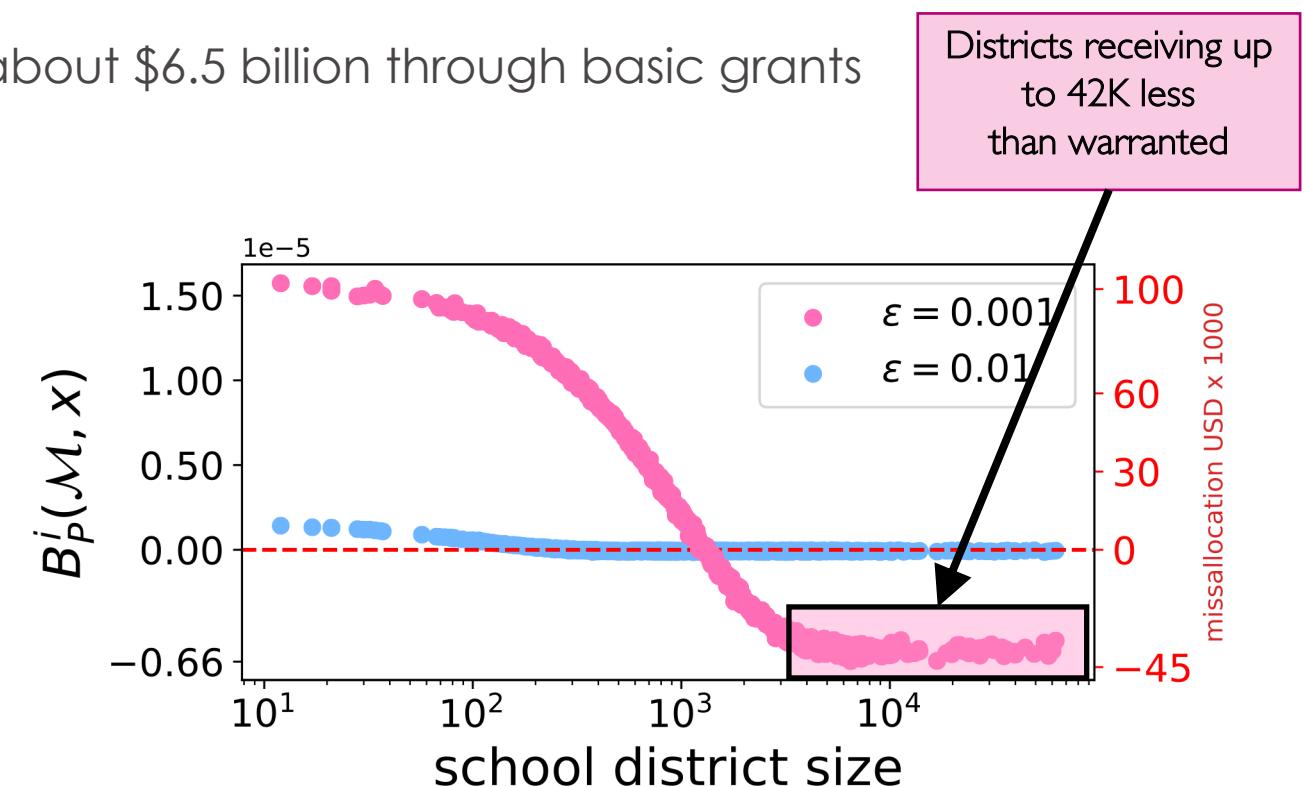
Disproportionate Effects in Title 1 Allotment

- Title 1 of the Elementary and Secondary Education Act is one of the largest U.S. program offering educational assistance to disadvantaged children
- In the fiscal year 2015 alone, it distributed about \$6.5 billion through basic grants
- Allotment:

count of children 5 to 17 in district i

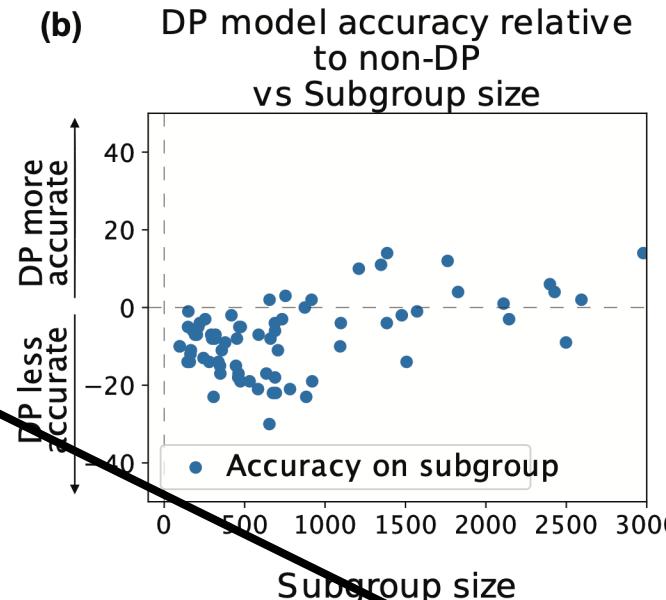
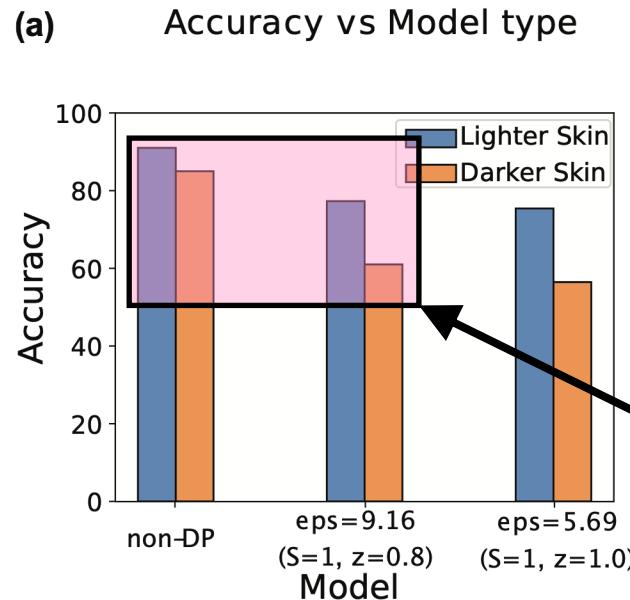
$$P_i^F(x) \stackrel{\text{def}}{=} \left(\frac{x_i \cdot a_i}{\sum_{i \in [n]} x_i \cdot a_i} \right)$$

student expenditures in district i

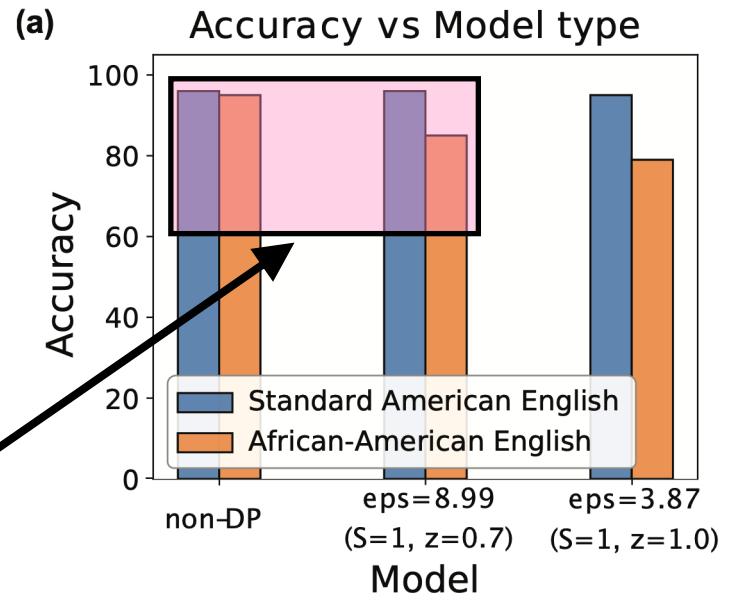


Disproportionate Effects in ML

Gender and age classification on facial images



Sentiment analysis of tweets

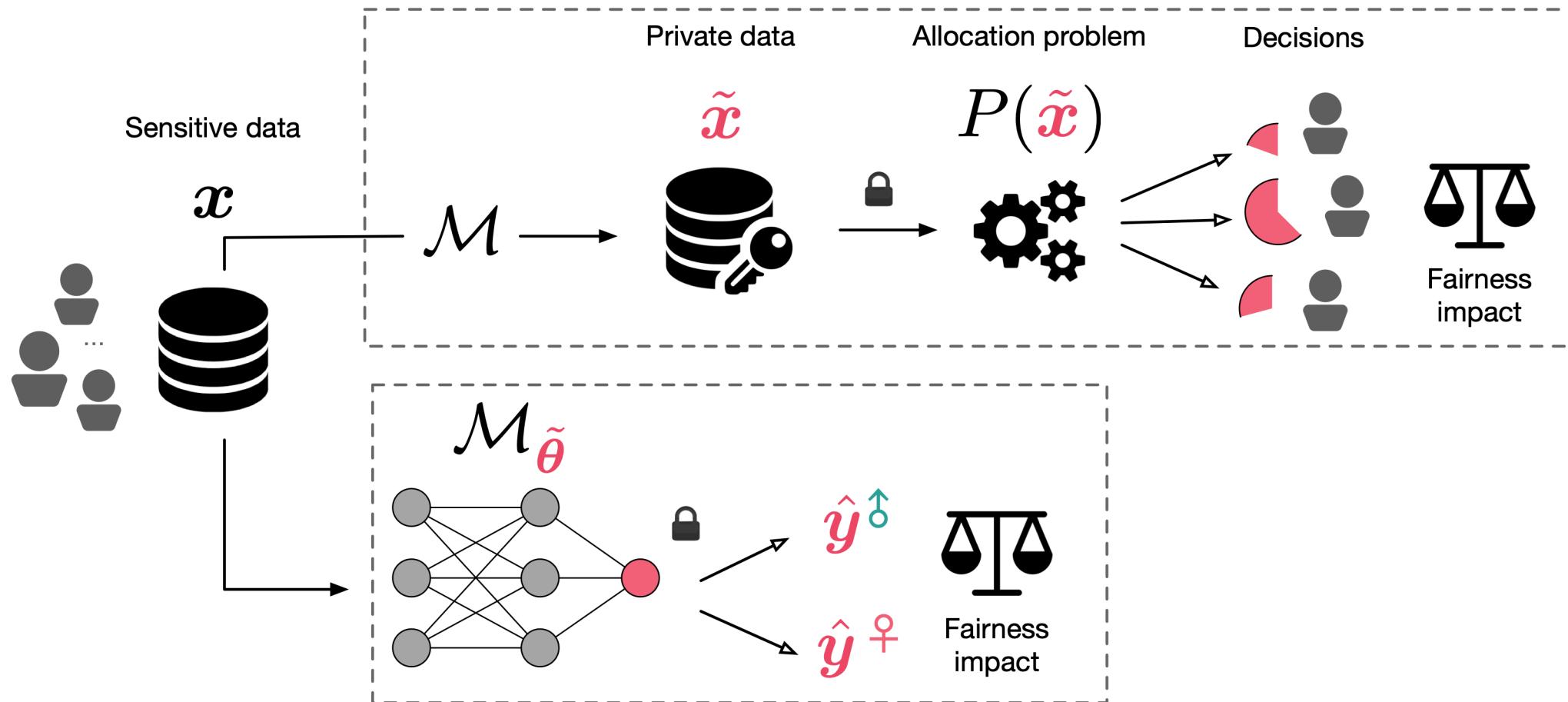


As the privacy increases the accuracy disparity of the learning task increases

Societal Impact

- The resulting outcomes can have significant societal and economic impacts on the involved individuals:
- **Classification errors** may penalize some groups over others in important determinations including criminal assessment, hiring or landing
- **Biased decisions** can result in disparities regarding the allocation of critical funds, benefits, and therapeutics
- While these observations are increasingly apparent, their causes have only recently started to receive attention.

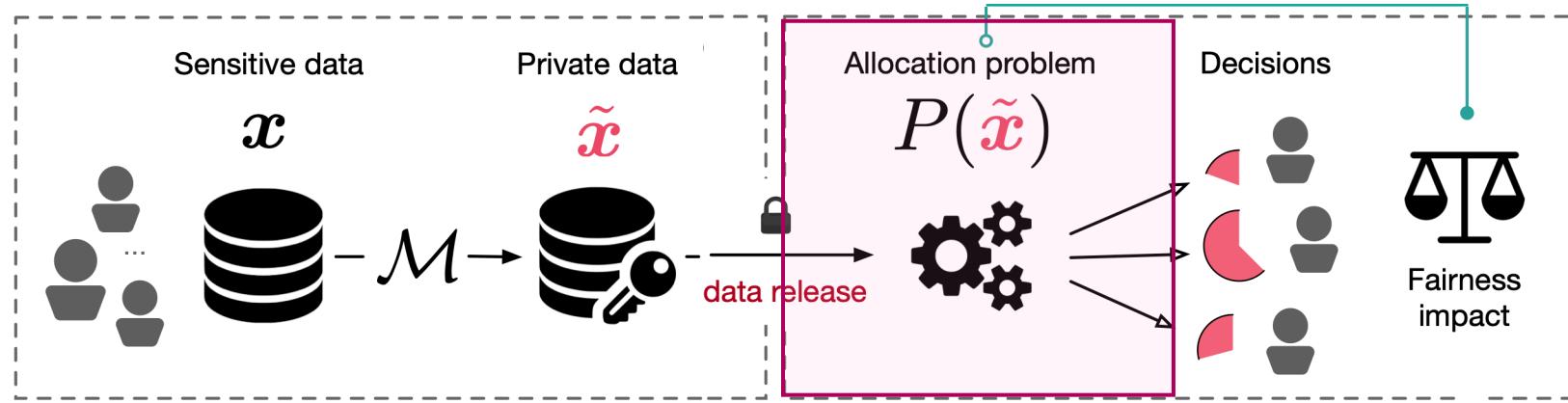
Setting and Outline of Part II



Why does disparity arise in decision making?

Census data-release perspective

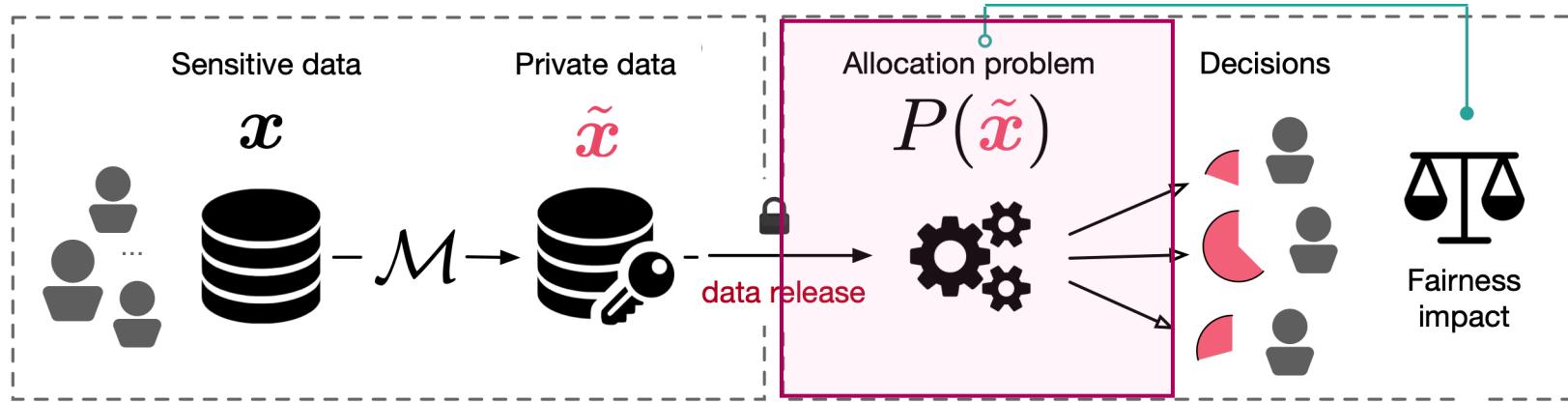
Fairness in DP Downstream Decisions



Given a dataset $x \in \mathcal{X} \subseteq \mathbb{R}^k$ of n entities, whose elements x_i describe some measurable quantity of entity $i \in [n]$ we consider two classes of decision problems:

1. **Allotment problems**: $P: \mathcal{X} \times [n] \rightarrow \mathbb{R}$,
which distributes a finite set of resources to some entity.
2. **Decision rules**: $P: \mathcal{X} \times [n] \rightarrow \{0,1\}$,
which determines whether an entity qualifies for some benefit.

Fairness in DP Downstream Decisions



Bias: $B_P^i(\mathcal{M}, x) = \mathbb{E}_{\tilde{x} \sim \mathcal{M}(x)} [P_i(\tilde{x})] - P_i(x).$

Definition (α -Fairness). A data-release mechanism \mathcal{M} is said α -fair w.r.t. a problem P if, for all datasets $x \in \mathcal{X}$ and all $i \in [n]$

$$\xi_B^i(P, \mathcal{M}, x) = \max_{j \in [n]} |B_P^i(\mathcal{M}, x) - B_P^j(\mathcal{M}, x)| \leq \alpha,$$

Fair Allotments

- First (surprising) result:

Even with **an unbiased DP mechanism**, the “shape” of the decision problem characterizes the unfairness of the outcomes.

Theorem 3. Let P be an allotment problem that is at least twice differentiable. A data-release mechanism \mathcal{M} is α -fair w.r.t. P , for some finite α , if for all datasets $x \in \mathcal{X}$ the entries of the Hessian HP_i of problem P_i are a constant function, that is, if there exists $c_{jl}^i \in \mathbb{R}$ ($i \in [n], j, l \in [k]$) such that,

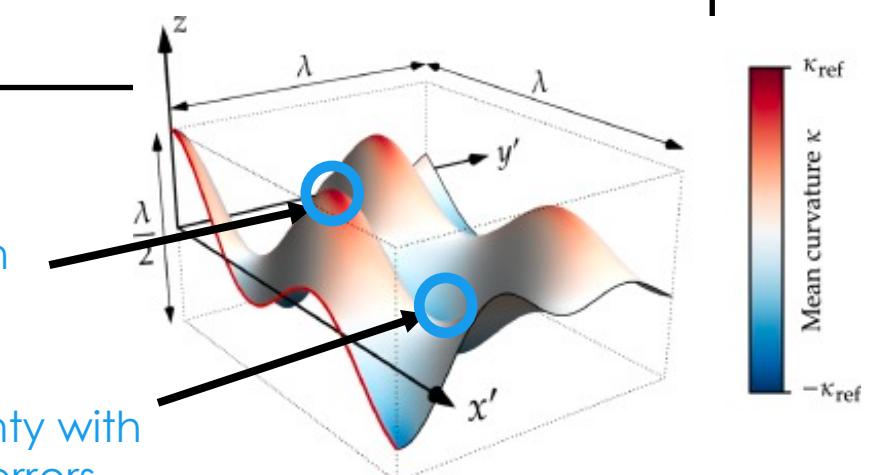
Hessian of problem P_i \rightarrow $(HP_i)_{j,l}(x) = c_{j,l}^i$ ($i \in [n] j, l \in [k]$).

Corollary 1 (informal). (Perfect)-fairness cannot be achieved if P is any non-convex function, as in the case of the allocations considered.

Adding Laplace noise to the inputs will necessarily introduce fairness issues, despite the noise being unbiased!

county with high errors

county with low errors



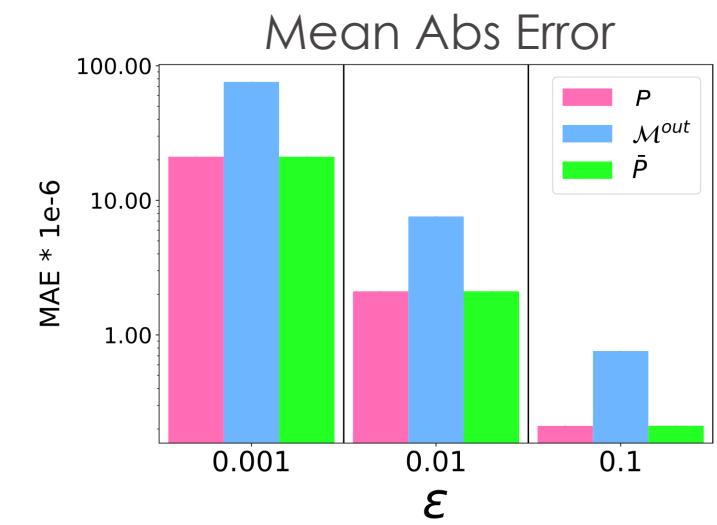
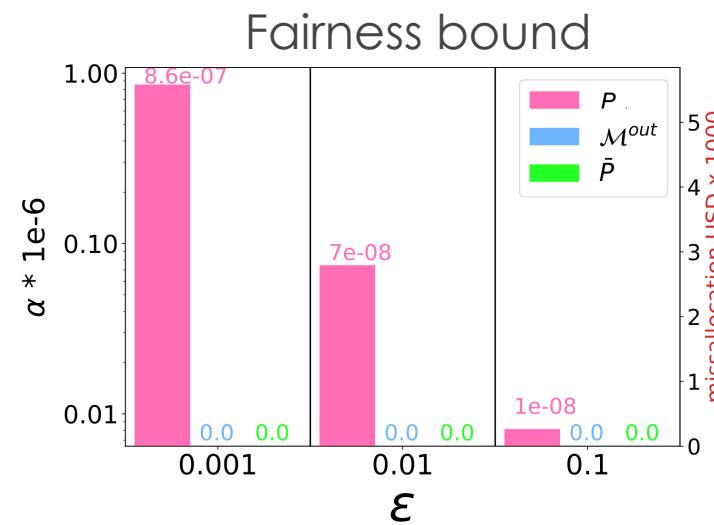
Fair Allotments (mitigation)

Corollary 2 (informal). If P is a linear function, then mechanism M is fair w.r.t. P

- **Note:** The observed issues are not data-driven, but problem-driven.
- **Linearizing the allotment problem** — General idea: Given a problem P_i derive a linear approximation \tilde{P}_i of P_i . For example, using a redundant data release:

$$P_i^F(x) \stackrel{\text{def}}{=} \left(\frac{x_i \cdot a_i}{\sum_{i \in [n]} x_i \cdot a_i} \right)$$

release its (noisy) version
as a constant



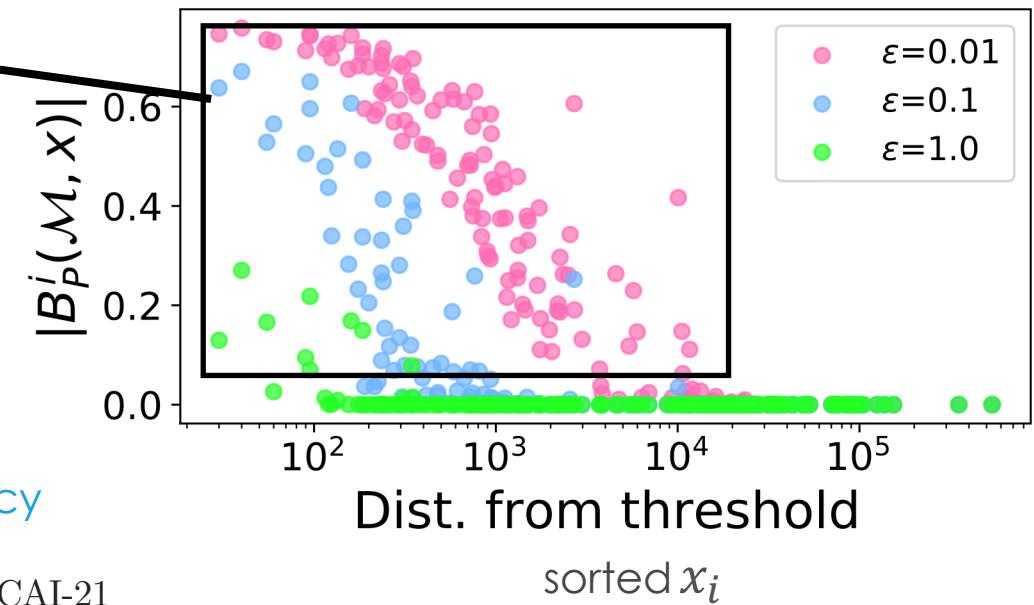
Disproportionate Effects in Minority Language Voting Rights

- The Voting Rights Act of 1965 provides a body of protections for racial and language minorities.
- Section 203 describes the conditions under which local jurisdictions must provide minority language voting assistance during an election.
- Jurisdiction i must provide language assistance (including voter registration, ballots, and instructions) iff decision rule $P(\mathbf{x})$ returns true with

$$P_i^M(\mathbf{x}) \stackrel{\text{def}}{=} \left(\frac{x_i^{sp}}{x_i^s} > 0.05 \vee x_i^{sp} > 10^4 \right) \wedge \frac{x_i^{spe}}{x_i^{sp}} > 0.0131.$$

+ $< 5^{\text{th}}$ grade education
 no. of ppl in i speaking minority language s
 + limited English proficiency

Misclassification implies potentially disenfranchising

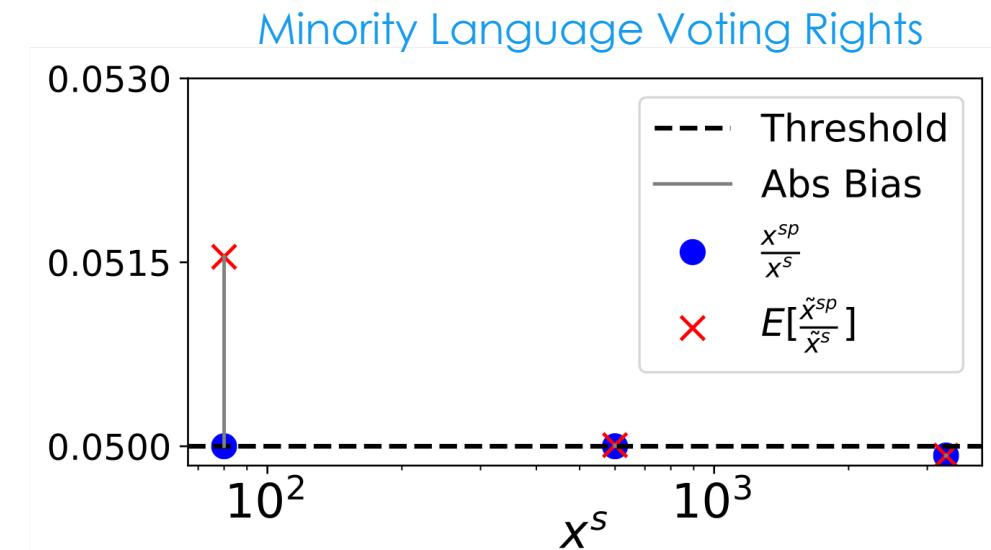


Fair Decision Rules

Ratio Functions

$$P_i^M(x) \stackrel{\text{def}}{=} \left(\frac{x_i^{sp}}{x_i^s} > 0.05 \vee x_i^{sp} > 10^4 \right) \wedge \frac{x_i^{spe}}{x_i^{sp}} > 0.0131.$$

- Loving county, TX, where $\frac{x^{sp}}{x^s} = 0.05$*
- Terrell county, TX, where $\frac{x^{sp}}{x^s} = 0.05$*
- Union county, NM, where $\frac{x^{sp}}{x^s} = 0.0484$*

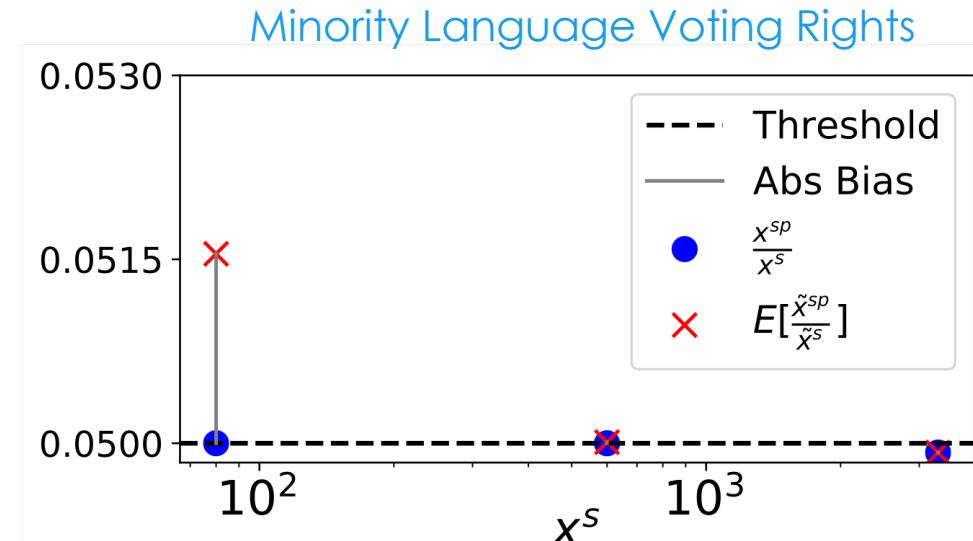


Fair Decision Rules

Ratio Functions

$$P_i^M(x) \stackrel{\text{def}}{=} \left(\frac{x_i^{sp}}{x_i^s} > 0.05 \vee x_i^{sp} > 10^4 \right) \wedge \frac{x_i^{spe}}{x_i^{sp}} > 0.0131.$$

- Loving county, TX, where $\frac{x^{sp}}{x^s} = \frac{4}{80} = 0.05$
- Terrell county, TX, where $\frac{x^{sp}}{x^s} = \frac{30}{600} = 0.05$
- Union county, NM, where $\frac{x^{sp}}{x^s} = \frac{160}{3305} = 0.0484$



Theorem (informal). The perturbation induced by the DP mechanism affects more the county with lower numerator / denominator.

Fair Decisions Rules (cont.)

- Second First (surprising) result:

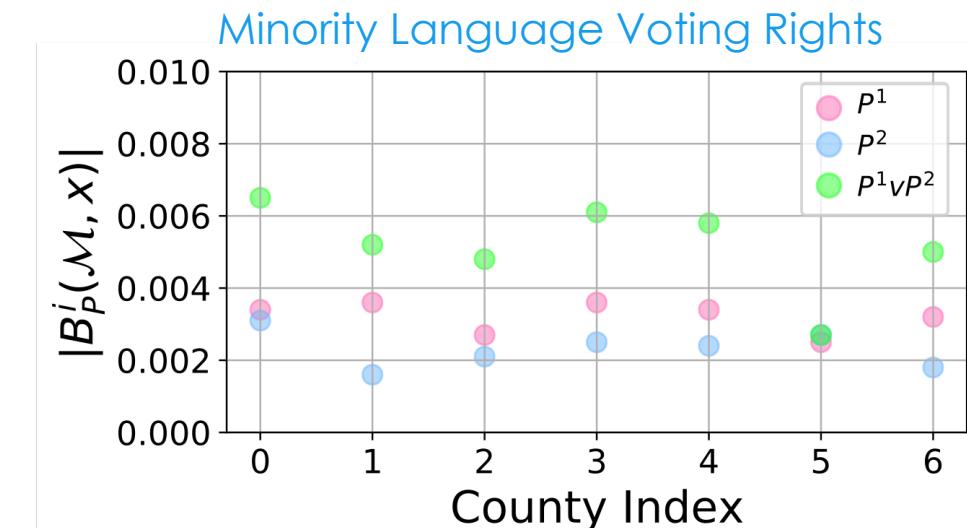
The unfairness induced by “composing” predicates is **larger than** that of the individual components.

$$P_i^M(x) \stackrel{\text{def}}{=} \left(\frac{x_i^{sp}}{x_i^s} > 0.05 \vee x_i^{sp} > 10^4 \right) \wedge \frac{x_i^{spe}}{x_i^{sp}} > 0.0131.$$

$$P^1(x^{sp}) = \mathbb{1}\{x^{sp} \geq 10^4\}$$

$$P^2(x^{sp}, x^{spe}) = \mathbb{1}\left\{ \frac{x^{spe}}{x^{sp}} > 0.0131 \right\}$$

Theorem (informal). The logical composition of two α_1 - and α_2 -fair mechanisms is α -fair with $\alpha > \max(\alpha_1, \alpha_2)$.



- Small bias when considered individually
- However, when combined using logical connector \wedge , the resulting absolute bias increases substantially (green circles).

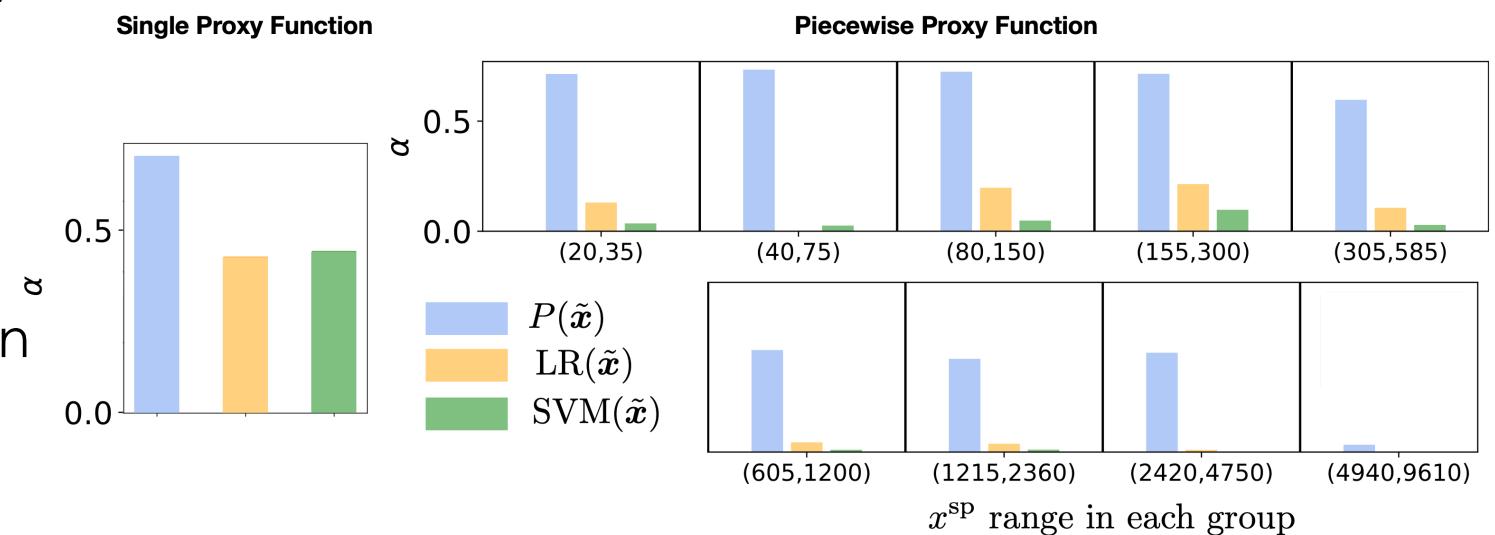
Fair Decisions (mitigation)

- Much more difficult scenario. But we could resort to the linear approximation trick again

- Partition dataset into groups x^{sp}
- Train subgroups using features \mathbf{x} using a linear classifier
- Use the parameters of the proxy linear model $LR(\mathbf{x})$ or $SVM(\mathbf{x})$ to make a decision i.e., to approximate P_i^M

$$P_i^M(\mathbf{x}) \stackrel{\text{def}}{=} \left(\frac{x_i^{sp}}{x_i^s} > 0.05 \vee x_i^{sp} > 10^4 \right) \wedge \frac{x_i^{spe}}{x_i^{sp}} > 0.0131.$$

Result summary: Fairness violation decreases substantially, albeit within each subgroup

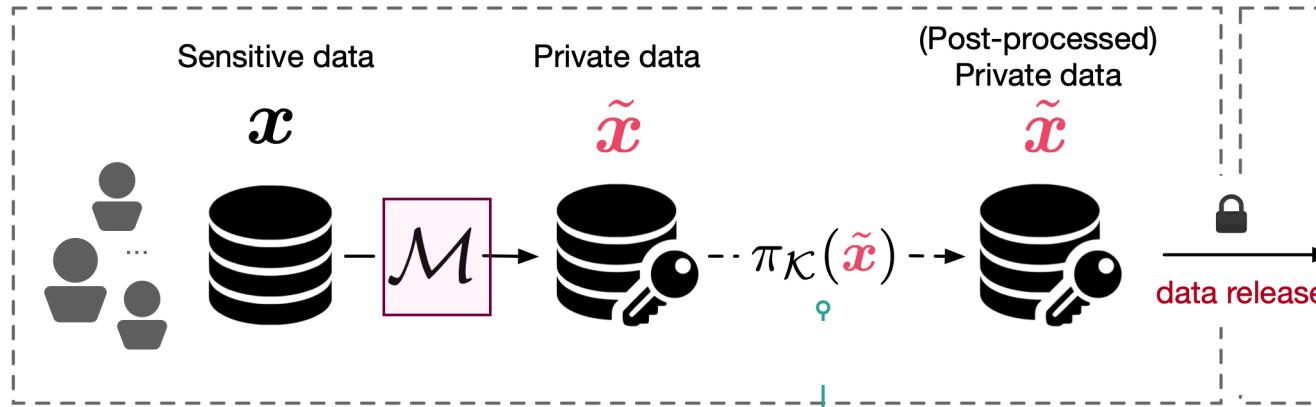


Why does disparity arise in decision making?

So far:

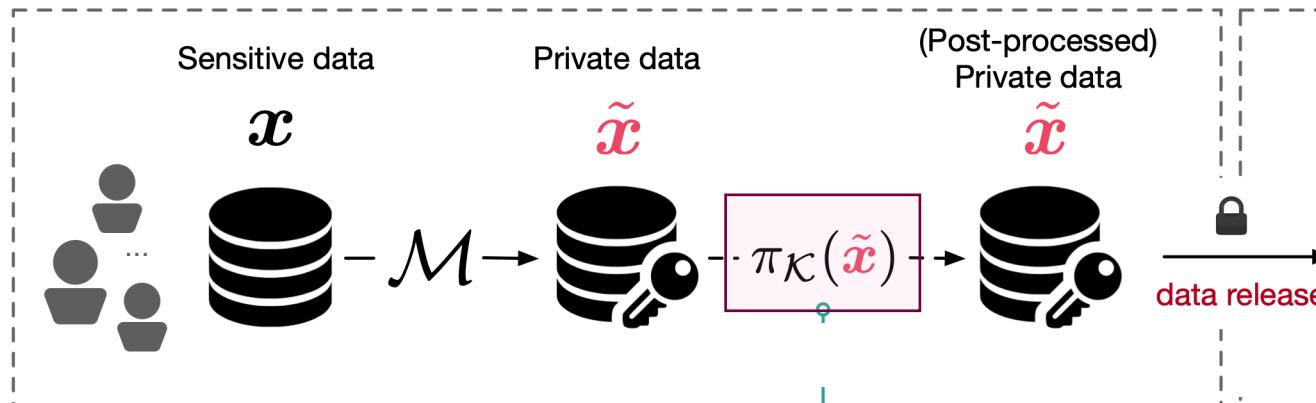
The “shape” of the downstream decision problem

Fairness in DP Downstream Decisions



I. Apply noise with appropriate parameter $\tilde{x} = x + \text{Noise}$

Fairness in DP Downstream Decisions



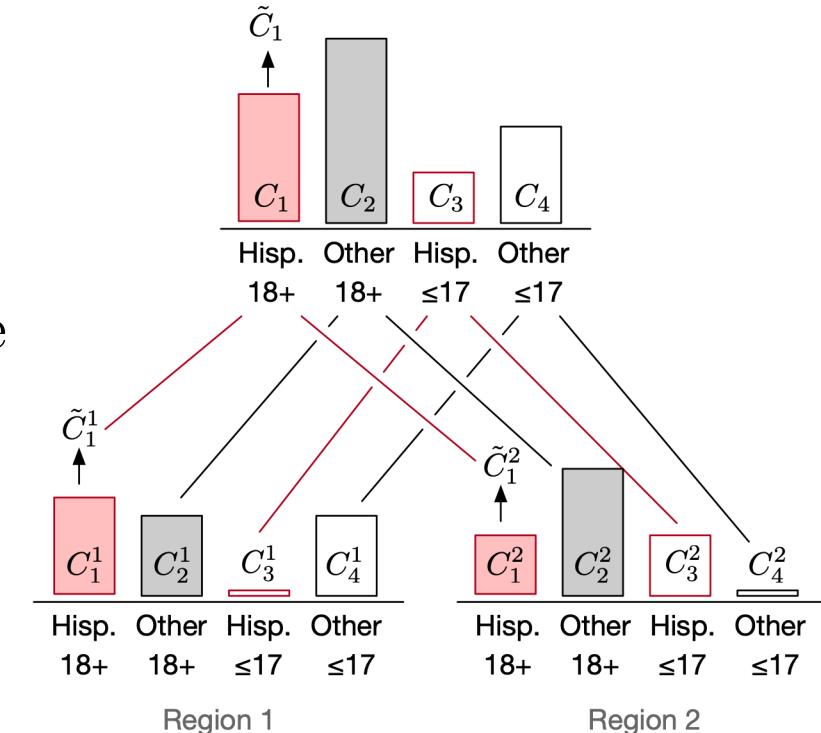
1. Apply noise with appropriate parameter $\tilde{x} = x + \text{Noise}$

2. Postprocess output \tilde{x} to enforce consistency

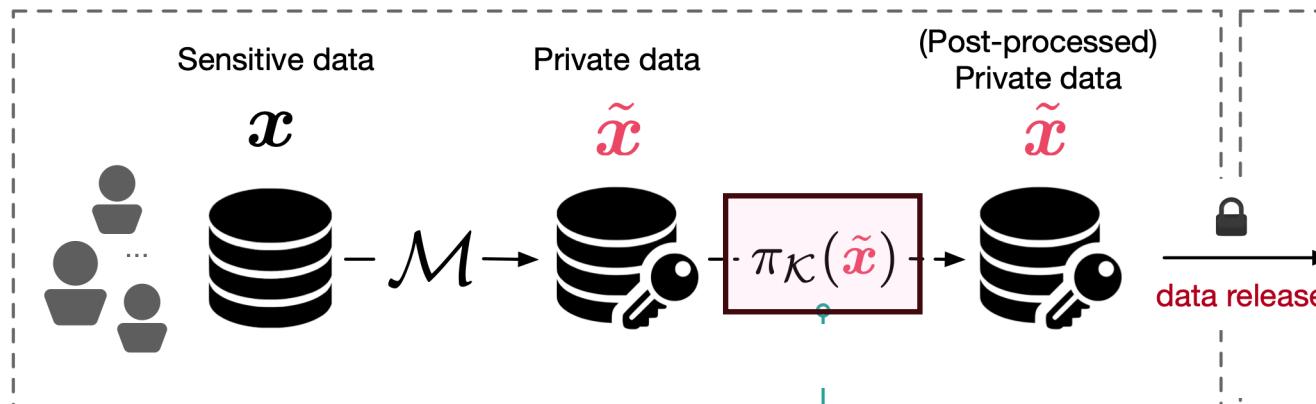
$$\pi_{S+}(\tilde{x}) := \arg \min_{v \in \mathcal{K}_{S+}} \|v - \tilde{x}\|_2,$$

with feasible region defined as

$$\mathcal{K}_{S+} = \left\{ v \mid \sum_{i=1}^n v_i = C, v \geq 0 \right\}.$$



Fairness in DP Downstream Decisions



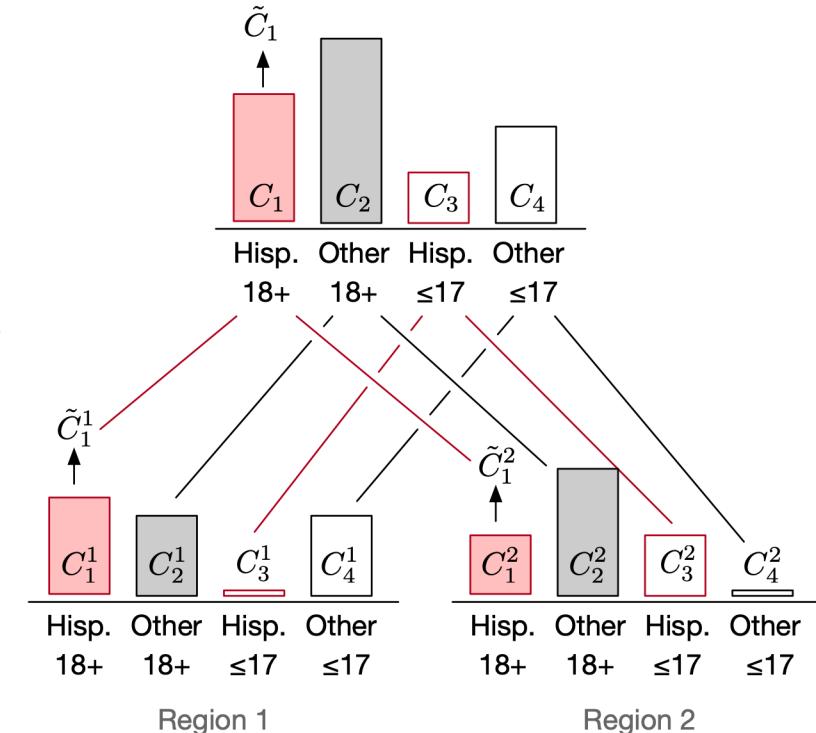
1. Apply noise with appropriate parameter $\tilde{x} = x + \text{Noise}$

2. Postprocess output \tilde{x} to enforce consistency

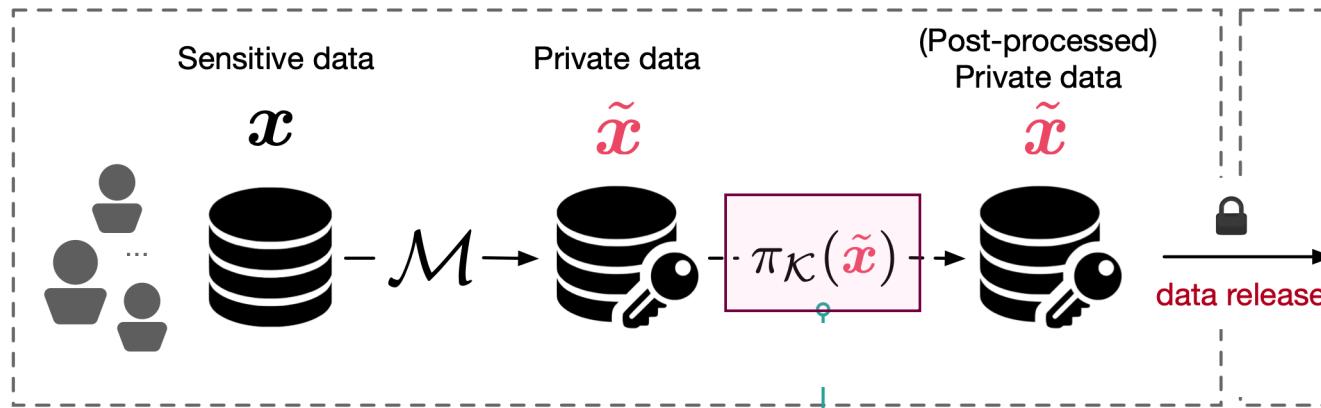
$$\pi_{S+}(\tilde{x}) := \arg \min_{\mathbf{v} \in \mathcal{K}_{S+}} \|\mathbf{v} - \tilde{x}\|_2,$$

with feasible region defined as

$$\mathcal{K}_{S+} = \left\{ \mathbf{v} \mid \sum_{i=1}^n v_i = C, \mathbf{v} \geq \mathbf{0} \right\}.$$



Fairness in DP Downstream Decisions



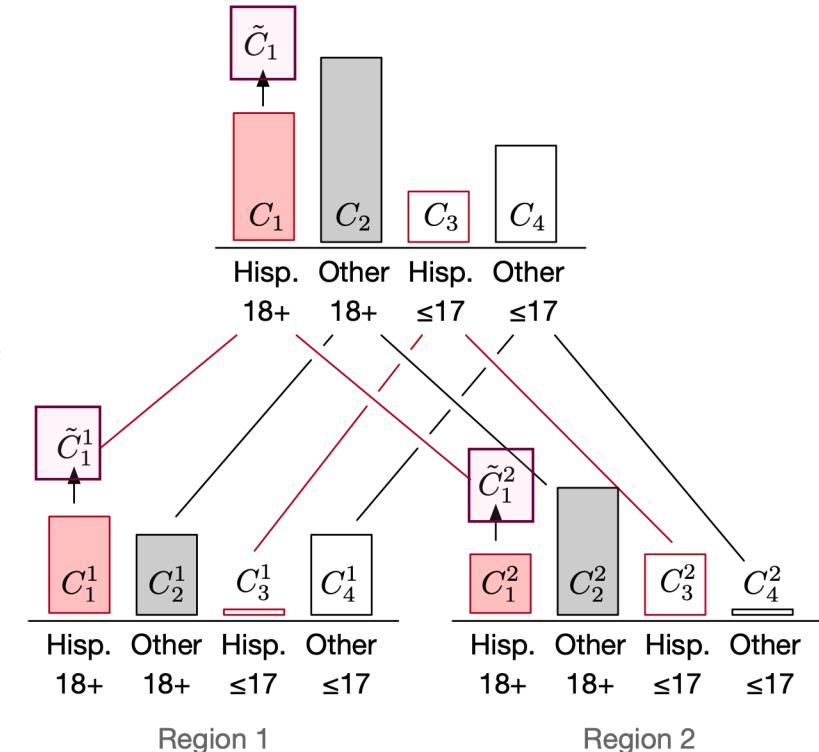
1. Apply noise with appropriate parameter $\tilde{\mathbf{x}} = \mathbf{x} + \text{Noise}$

2. Postprocess output $\tilde{\mathbf{x}}$ to enforce consistency

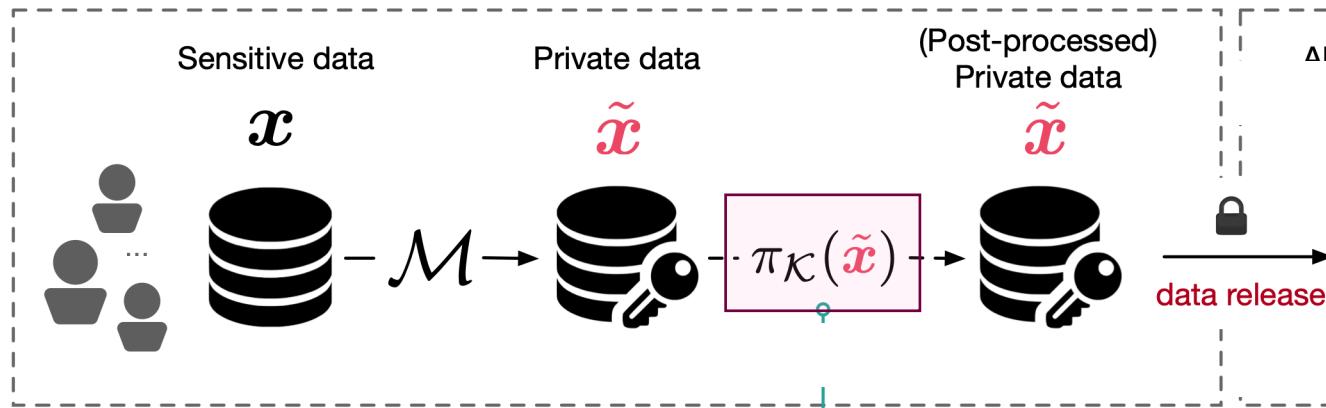
$$\pi_{S+}(\tilde{\mathbf{x}}) := \arg \min_{\mathbf{v} \in \mathcal{K}_{S+}} \|\mathbf{v} - \tilde{\mathbf{x}}\|_2,$$

with feasible region defined as

$$\mathcal{K}_{S+} = \left\{ \mathbf{v} \mid \sum_{i=1}^n v_i = C, \mathbf{v} \geq \mathbf{0} \right\}.$$



Fairness in DP Downstream Decisions



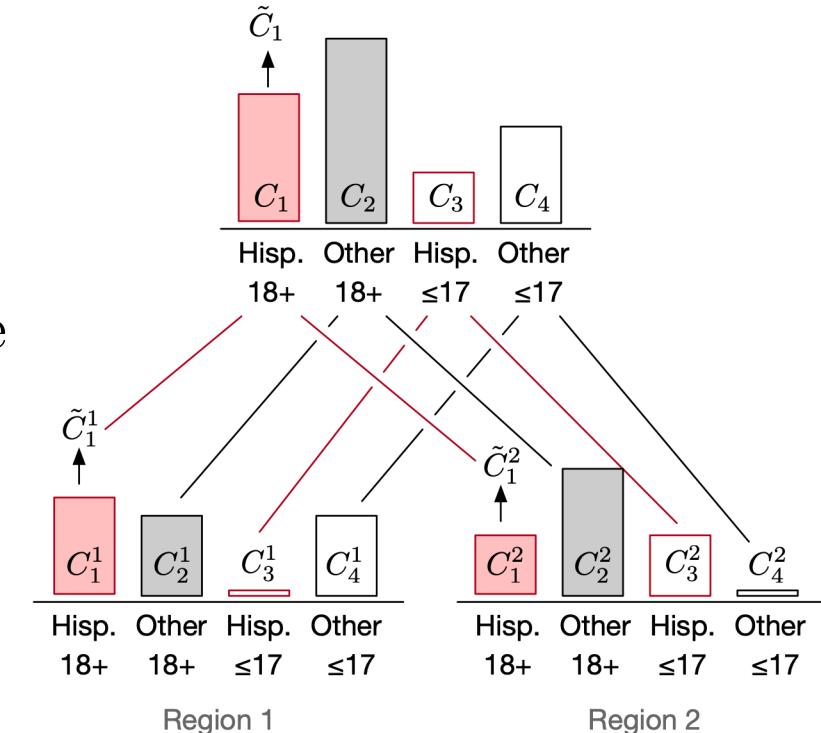
I. Apply noise with appropriate parameter $\tilde{\mathbf{x}} = \mathbf{x} + \text{Noise}$

2. Postprocess output $\tilde{\mathbf{x}}$ to enforce consistency

$$\pi_{S+}(\tilde{\mathbf{x}}) := \arg \min_{\mathbf{v} \in \mathcal{K}_{S+}} \|\mathbf{v} - \tilde{\mathbf{x}}\|_2,$$

with feasible region defined as

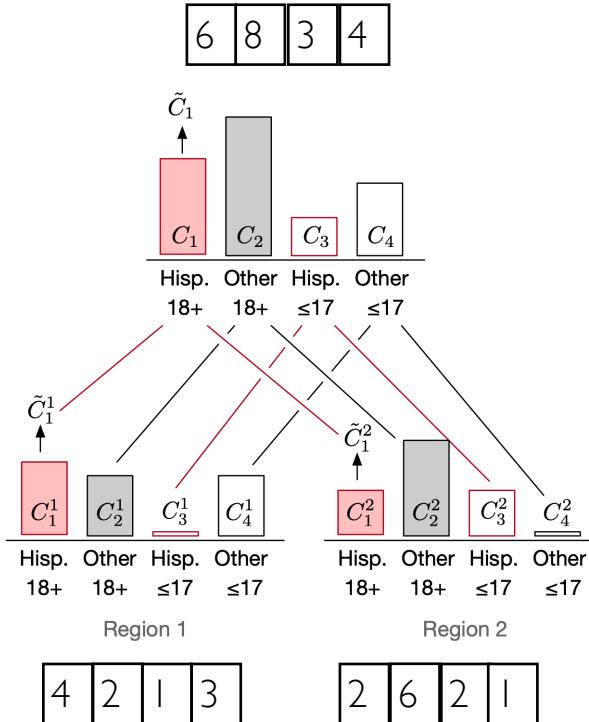
$$\mathcal{K}_{S+} = \left\{ \mathbf{v} \mid \sum_{i=1}^n v_i = C, \boxed{\mathbf{v} \geq \mathbf{0}} \right\}.$$



Satisfies DP due to post-processing immunity

DP Post-Processing

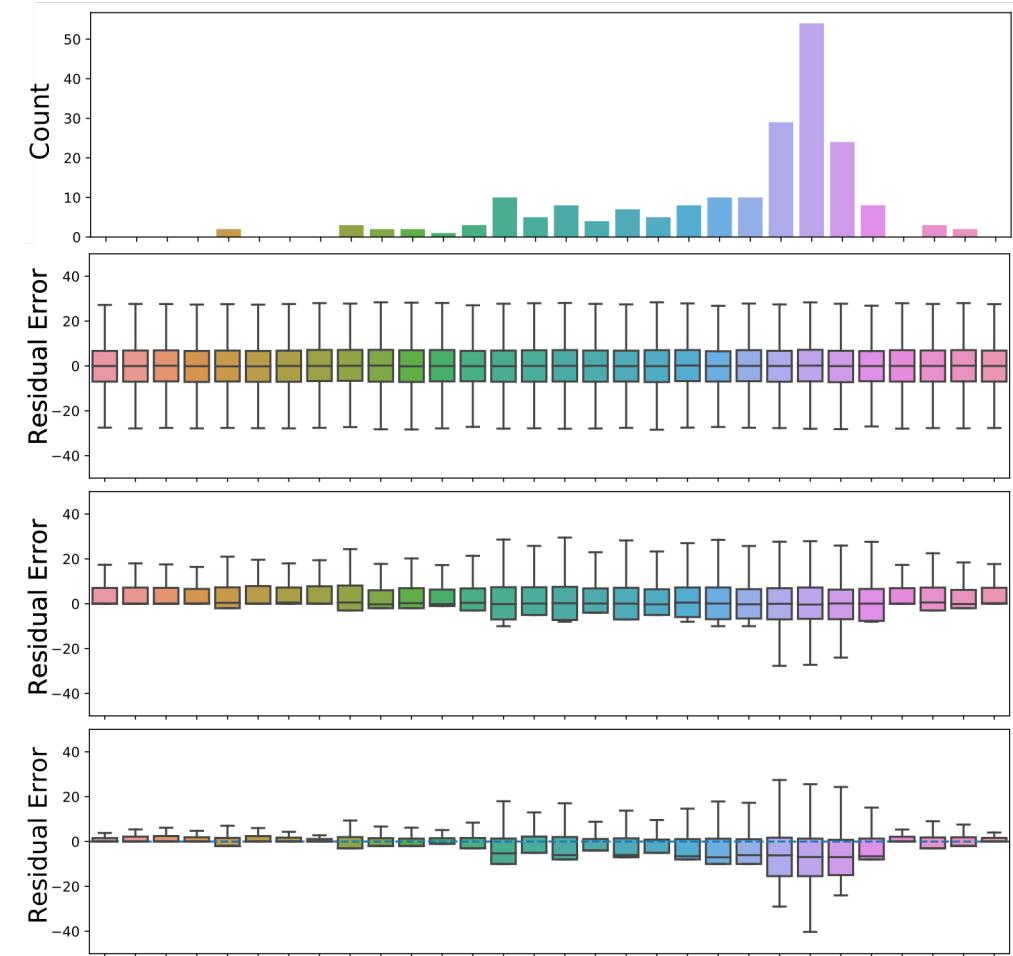
Post-processing



$$\pi_{\geq 0} := \operatorname{argmin}_{v \geq 0} \|v - \tilde{x}\|_2$$

$$\pi_{\mathcal{K}_S} := \operatorname{argmin}_{v \in \mathcal{K}_S} \|v - \tilde{x}\|_2, \quad \mathcal{K}_S = \{v \in \mathbb{R}^n \mid \sum_i v_i = \tilde{s}, v_i \geq 0\},$$

Laplace
mechanism



DP Post-Processing

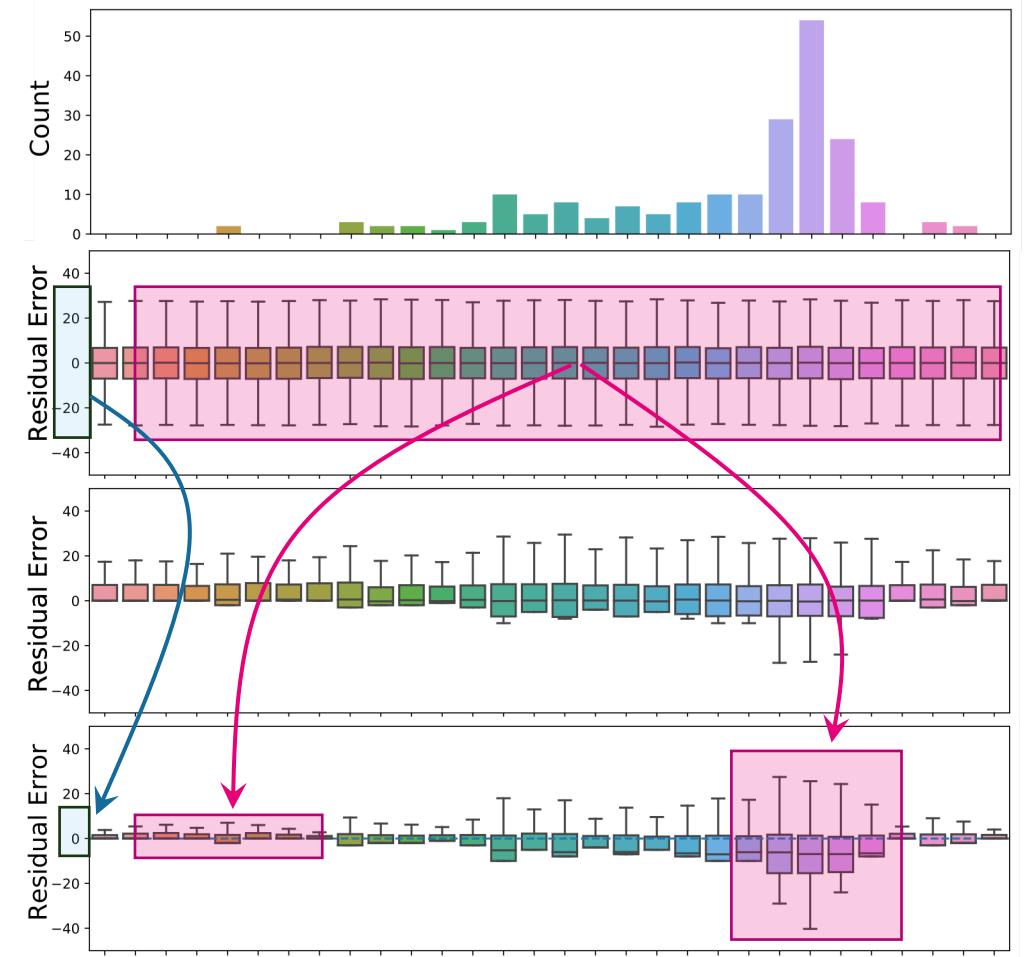
Third result: Observe that post-processing reduce the errors;

However, it increases unfairness!

Laplace mechanism

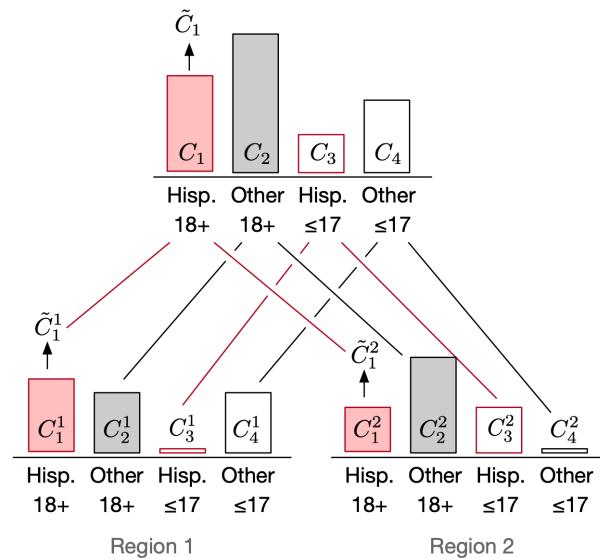
$$\pi_{\geq 0} := \operatorname{argmin}_{v \geq 0} \|v - \tilde{x}\|_2$$

$$\pi_{\mathcal{K}_S} := \operatorname{argmin}_{v \in \mathcal{K}_S} \|v - \tilde{x}\|_2, \quad \mathcal{K}_S = \{v \in \mathbb{R}^n \mid \sum_i v_i = \tilde{S}, v_i \geq 0\},$$



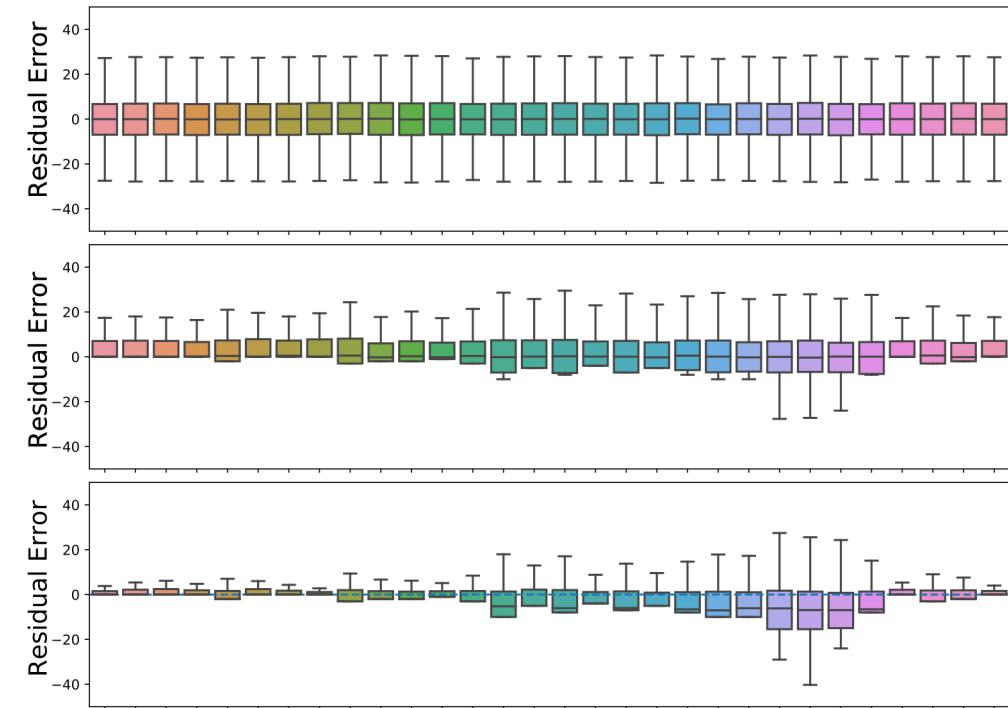
DP Post-Processing

Why is this happening? The bias comes from the **non-negativity constraints!**



$$\pi_{\geq 0} := \operatorname{argmin}_{v \geq 0} \|v - \tilde{x}\|_2$$

$$\pi_{\mathcal{K}_S} := \operatorname{argmin}_{v \in \mathcal{K}_S} \|v - \tilde{x}\|_2, \quad \mathcal{K}_S = \{v \in \mathbb{R}^n \mid \sum_i v_i = \tilde{S}, v_i \geq 0\},$$

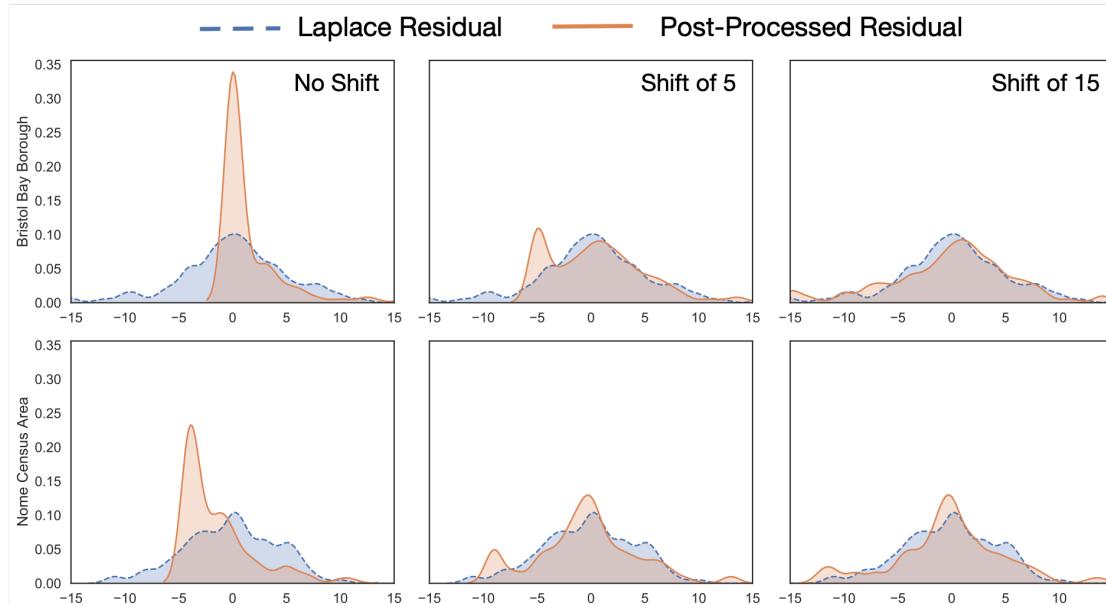


Quantifying Bias in DP Post-Processing

Theorem 6. Suppose that the noisy data \tilde{x} is the output of the Laplace mechanism with scale λ . The bias of the post-processed solution $\pi_{\mathcal{K}^+}$ of program (L^+) is bounded, in l_∞ norm, by

$$\|B_{L^+}(\mathcal{M}, \mathbf{x})\|_\infty = \left\| \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{M}(\mathbf{x})} [\pi_{L^+}(\tilde{\mathbf{x}}) - \mathbf{x}] \right\|_\infty \leq C' \cdot \exp\left(\frac{-r_m}{\lambda}\right) \cdot \sum_{i=0}^{n-1} \frac{(r_m)^i}{i! \cdot \lambda^i},$$

where C' represents the value $\sup_{v \in \mathcal{K}^+} \|v - \mathbf{x}\|_\infty$, which is finite due to the boundedness of the feasible region \mathcal{K}^+ .

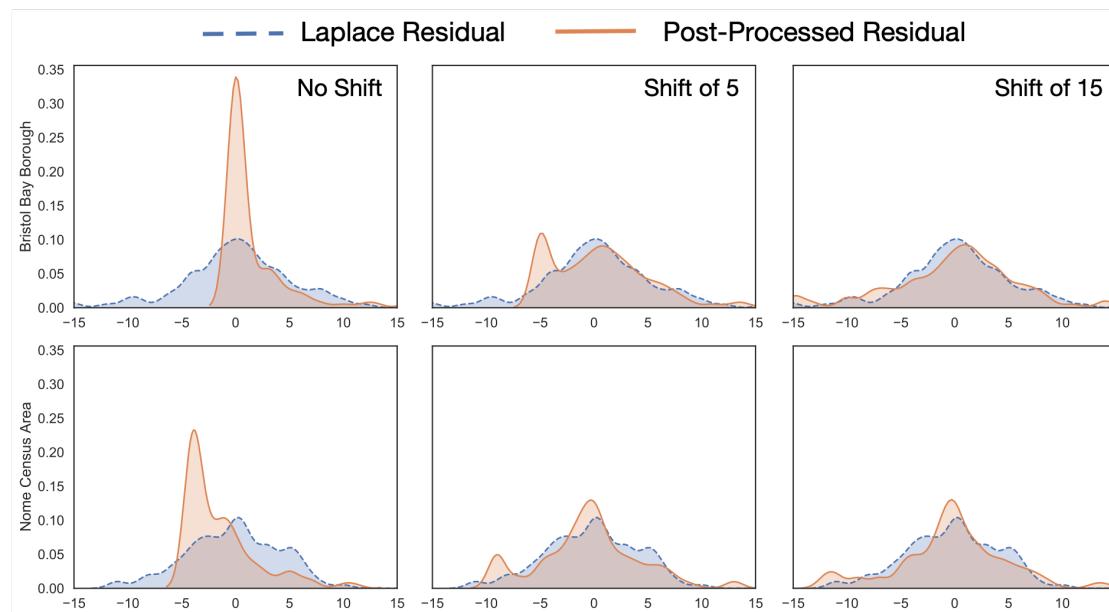


There is a L_1 -ball of radius $r_m = \min_i x_i$ and centered in \mathbf{x} which is a feasible subspace where there is no bias

Shifting increases the value of r_m and the bias progressively disappear

Quantifying Bias in DP Post-Processing

- What does this means practically? Post-processing reduces the variance of the noise differently in different “regions”. Regions with many subregions (e.g., counties, census blocks, etc.) will have more variance than regions with few subregions.
- It creates situations where counties will be treated fundamentally differently in decision processes.



Aggregating the counts for **Variance**

| | |
|---|--------|
| Arizona (pop: 2.37ML in 15 counties) | 186.67 |
| Texas (pop: 8.89ML in 254 counties) | 200.01 |

~6.5% difference
which may affect allocations!

Possible Mitigating Solution

Definition 4 (Projection onto Simplex Mechanism (PoS)).

The projection onto simplex mechanism *outputs the allocation as follows.*

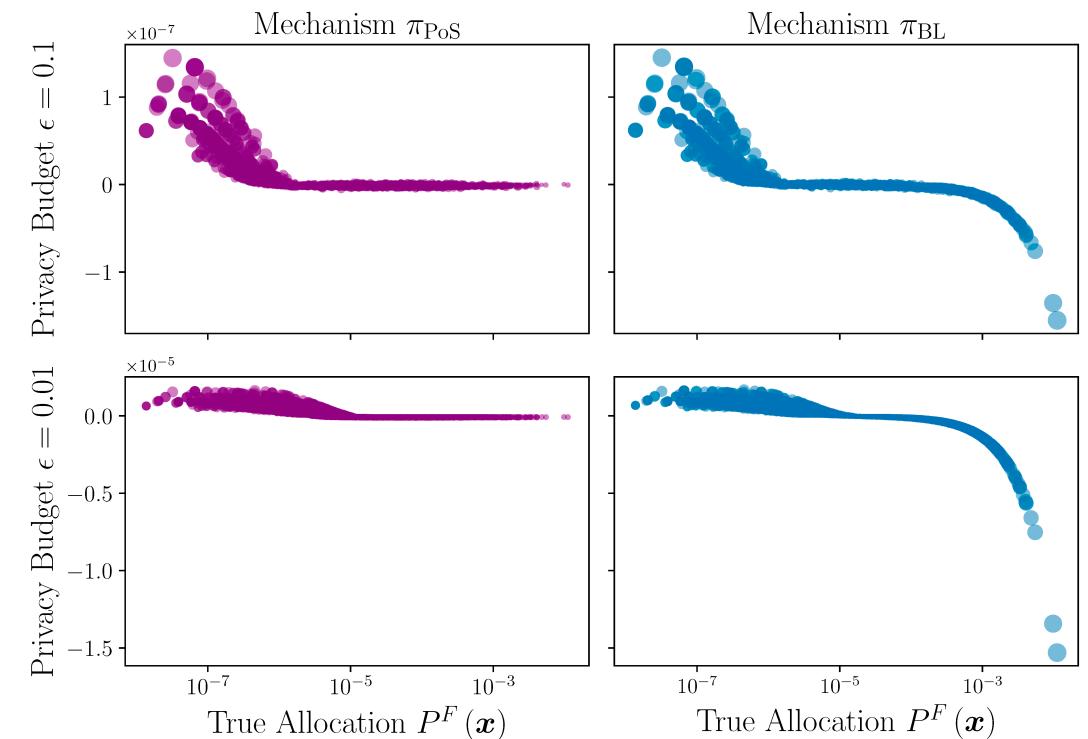
$$\pi_{\text{PoS}}(\tilde{\mathbf{x}}) := \arg \min_{\mathbf{v} \in \Delta^n} \|\mathbf{v} - P^F(\tilde{\mathbf{x}})\|_2 \quad (P_{\text{PoS}})$$

Theorem (informal). For any DP dataset $\tilde{\mathbf{x}}$ the PoS mechanism generates the unique optimal solution to program

$$\pi_\alpha^*(\tilde{\mathbf{x}}) := \arg \min_{\mathbf{v} \in \Delta_n} \|\mathbf{v} - P^F(\tilde{\mathbf{x}})\|_{\mathbb{W}} \quad (P_\alpha)$$

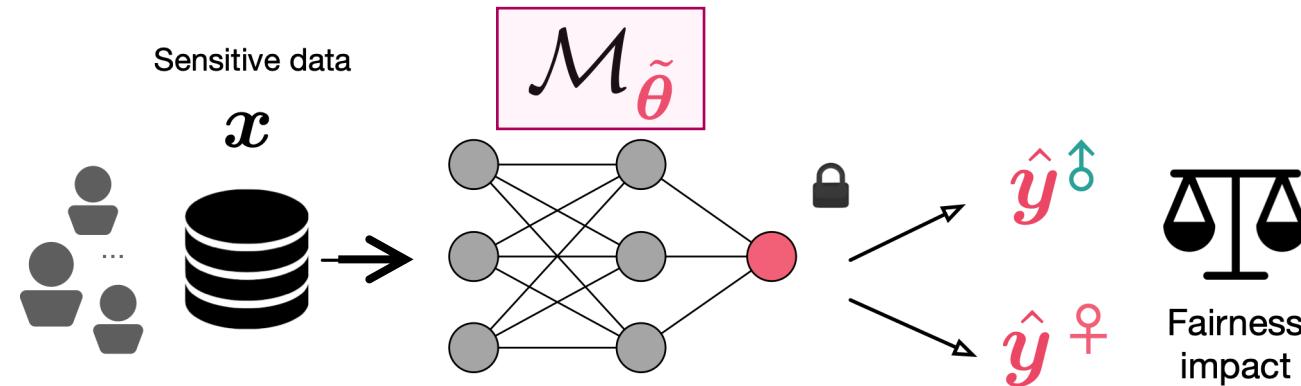
which closely approximate the optimal post-processing mechanism

$$\pi^* := \min_{\pi \in \Pi_{\Delta_n}} \|\mathbb{E}_{\tilde{\mathbf{x}}} [\pi(\tilde{\mathbf{x}}) - P^F(\mathbf{x})]\|_{\mathbb{W}}, \quad (4)$$



Why disparity arise in learning tasks?

Fairness in DP Learning Tasks



Given a dataset consisting of data points (X_i, A_i, Y_i) the goal is to learn a classifier f_θ that guarantees privacy of the individual data points and the learning task minimizes

$$\min_{\theta} \mathcal{L}(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(X_i), Y_i).$$

Fairness focuses on the notion of excessive risk: $R(\theta, D) = \mathbb{E}_{\tilde{\theta}} [\mathcal{L}(\tilde{\theta}; D)] - \mathcal{L}(\theta^*; D)$, and is measured with respect to the **excessive risk gap**

$$\xi_a = |R_a(\theta) - R(\theta)|.$$

↑ group-level ER ↑ Population-level ER

Warm up: Output Perturbation

Adds Gaussian noise $\mathcal{N}(0, \Delta_\ell^2 \sigma^2)$ to the optimal model parameters θ^* .

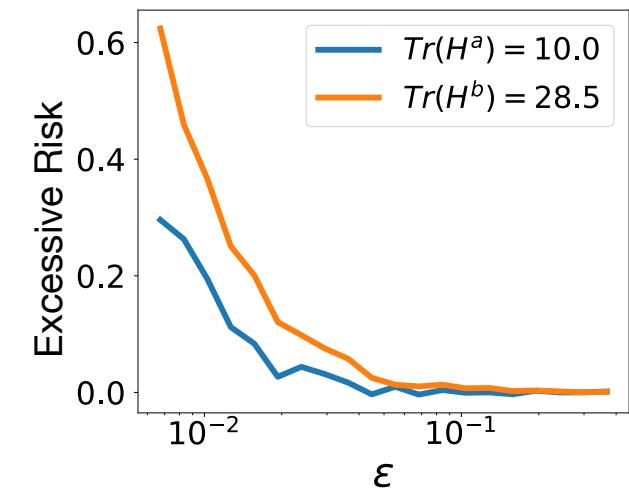
Theorem 1. Let ℓ be a twice differentiable and convex loss function and consider the output perturbation mechanism described above. Then, the excessive risk gap for group $a \in \mathcal{A}$ is approximated by:

$$\xi_a \approx \frac{1}{2} \Delta_\ell^2 \sigma^2 \left| \text{Tr}(\mathbf{H}_\ell^a) - \text{Tr}(\mathbf{H}_\ell) \right|, \quad (3)$$

where $\mathbf{H}_\ell^a = \nabla_{\theta^*}^2 \sum_{(X, A, Y) \in D_a} \ell(f_{\theta^*}(X), Y)$ is the Hessian matrix of the loss function, at the optimal parameters vector θ^* , computed using the group data D_a , \mathbf{H}_ℓ is the analogous Hessian computed using the population data D , and $\text{Tr}(\cdot)$ denotes the trace of a matrix.

Corollary 1. Consider the ERM problem for a linear model $f_\theta(X) = \theta^T X$, with L_2 loss function. Then, output perturbation does not guarantee pure fairness.

$$\text{Tr}(\mathbf{H}_\ell^a) = \mathbb{E}_{X \sim D_a} \text{Tr}(XX^T) = \mathbb{E}_{X \sim D_a} \|X\|^2$$



Warm up: Output Perturbation

Adds Gaussian noise $\mathcal{N}(0, \Delta_\ell^2 \sigma^2)$ to the optimal model parameters θ^* .

Theorem 1. Let ℓ be a twice differentiable and convex loss function and consider the output perturbation mechanism described above. Then, the excessive risk gap for group $a \in \mathcal{A}$ is approximated by:

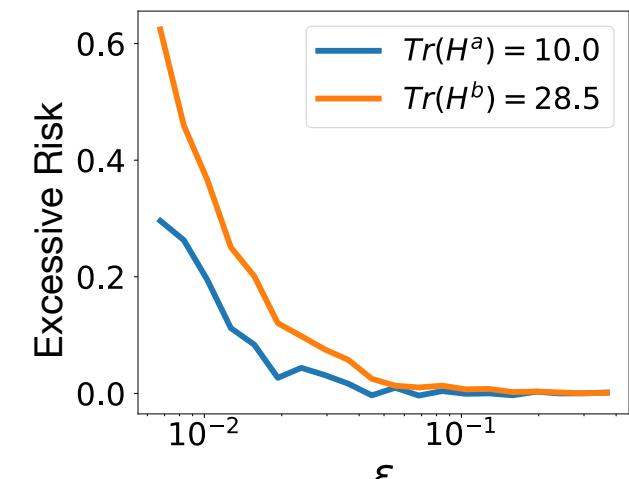
$$\xi_a \approx \frac{1}{2} \Delta_\ell^2 \sigma^2 \left| \text{Tr}(\mathbf{H}_\ell^a) - \text{Tr}(\mathbf{H}_\ell) \right|, \quad (3)$$

where $\mathbf{H}_\ell^a = \nabla_{\theta^*}^2 \sum_{(X, A, Y) \in D_a} \ell(f_{\theta^*}(X), Y)$ is the Hessian matrix of the loss function, at the optimal parameters vector θ^* , computed using the group data D_a , \mathbf{H}_ℓ is the analogous Hessian computed using the population data D , and $\text{Tr}(\cdot)$ denotes the trace of a matrix.

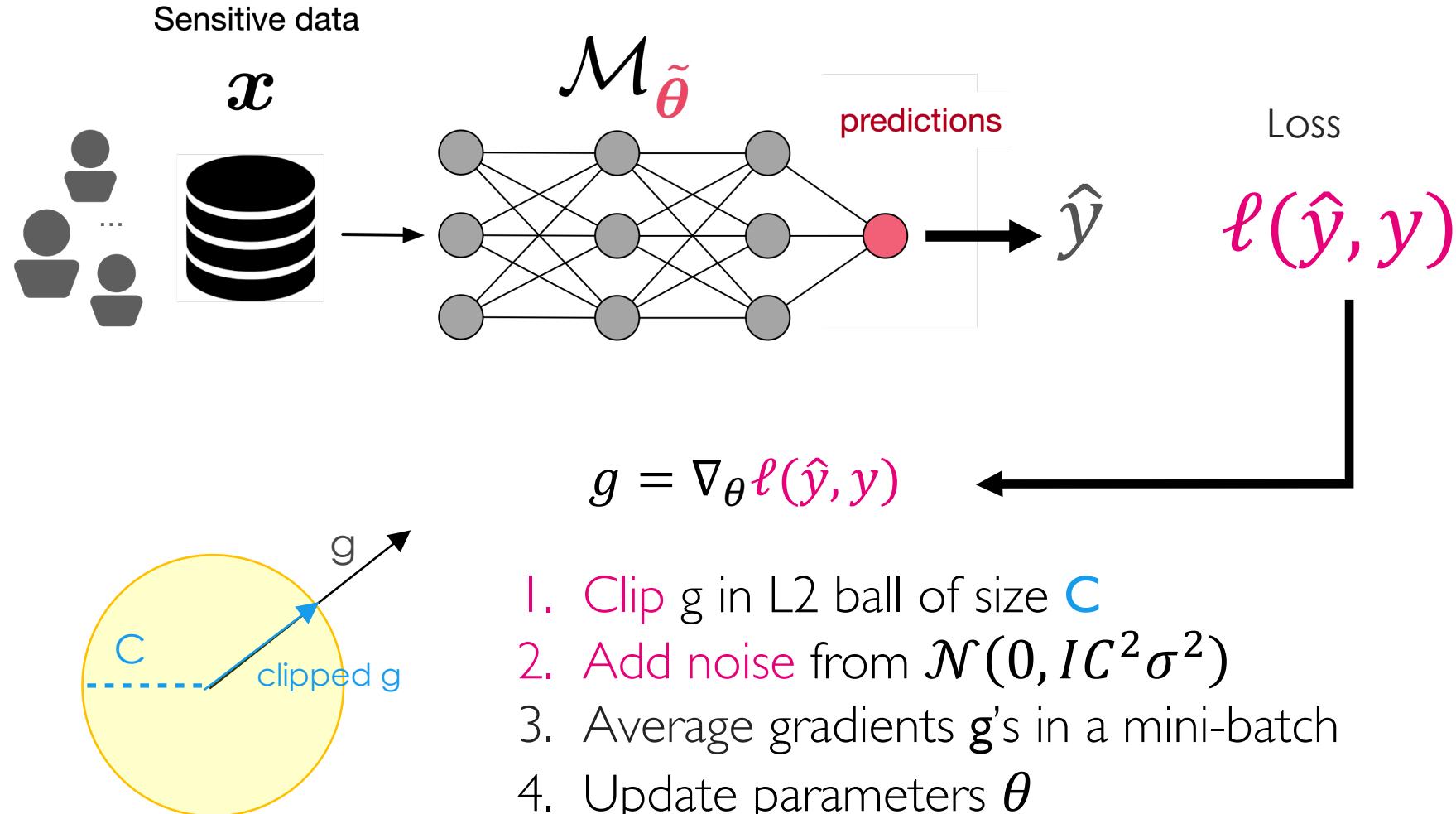
Corollary 2. For groups a and b, if their average group norms

$$\mathbb{E}_{X_a \sim D_a} \|X_a\| = \mathbb{E}_{X_b \sim D_b} \|X_b\|$$

have identical values, then output perturbation with L_2 loss function achieves pure fairness.



DP-Stochastic Gradient Descent



Fairness in DP-SGD

Theorem 2. Consider the ERM problem (L) with loss ℓ twice differentiable w.r.t. the model parameters. The expected loss $\mathbb{E}[\mathcal{L}(\theta_{t+1}; D_a)]$ of group $a \in \mathcal{A}$ at iteration $t+1$, is approximated as:

$$\begin{aligned} \mathbb{E}[\mathcal{L}(\theta_{t+1}; D_a)] &= \underbrace{\mathcal{L}(\theta_t; D_a) - \eta \langle g_{D_a}, g_D \rangle + \frac{\eta^2}{2} \mathbb{E}[g_B^T H_\ell^a g_B]}_{\text{non-private term}} \quad (4) \\ &\quad + \underbrace{\eta (\langle g_{D_a}, g_D \rangle - \langle g_{D_a}, \bar{g}_D \rangle) + \frac{\eta^2}{2} (\mathbb{E}[\bar{g}_B^T H_\ell^a \bar{g}_B] - \mathbb{E}[g_B^T H_\ell^a g_B])}_{\text{private term due to clipping}} \quad (R_a^{\text{clip}}) \\ &\quad + \underbrace{\frac{\eta^2}{2} \text{Tr}(H_\ell^a) C^2 \sigma^2}_{\text{private term due to noise}} \quad (R_a^{\text{noise}}) \\ &\quad + O(\|\theta_{t+1} - \theta_t\|^3), \end{aligned}$$

where the expectation is taken over the randomness of the private noise and the mini-batch selection, and the terms g_Z and \bar{g}_Z denote, respectively, the average non-private and private gradients over subset Z of D at iteration t (the iteration number is dropped for ease of notation).

Why clipping causes unfairness in DP-SGD

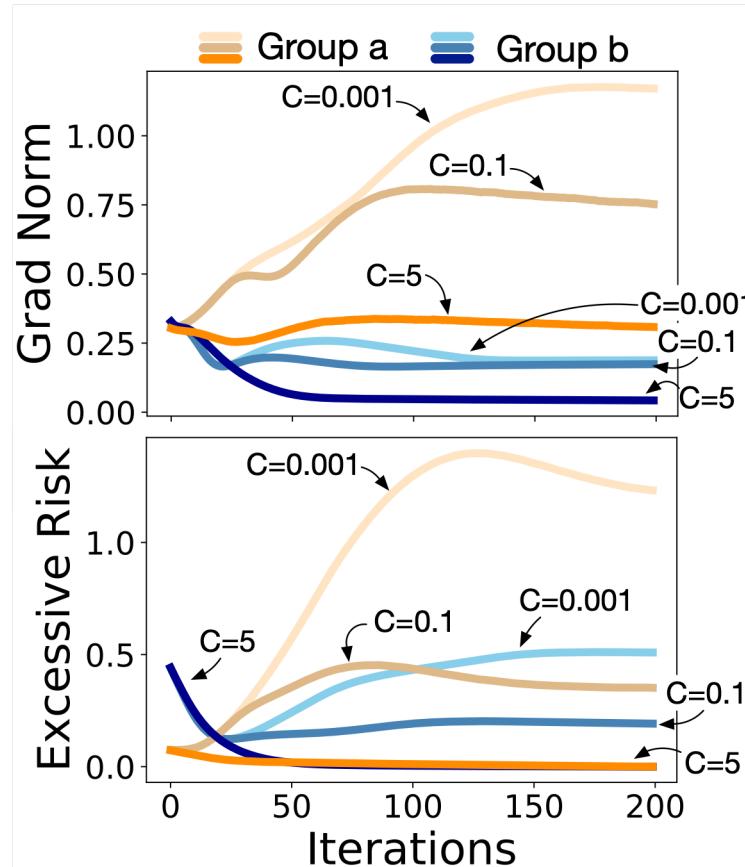
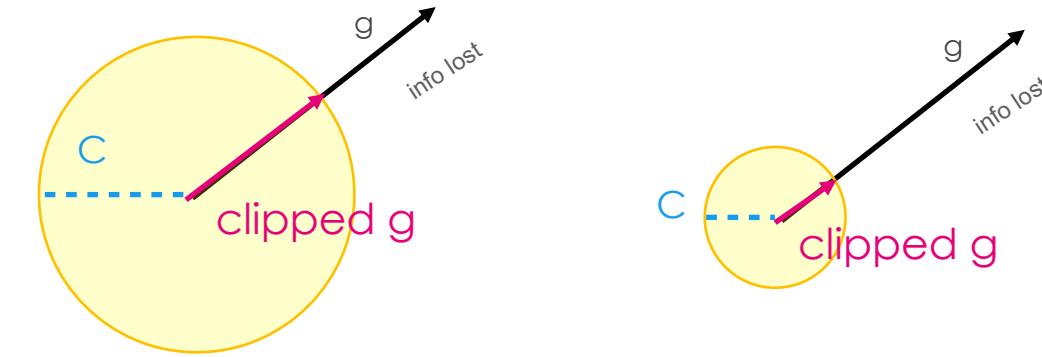


Figure 3: Impact of gradient clipping on gradient norms for different clipping bounds. Bank dataset.

When **clipping**, the smaller C, the higher
Is the information loss of the average gradients
that are backpropagated.



Theorem 3. Let $p_z = |D_z|/|D|$ be the fraction of training samples in group $z \in \mathcal{A}$. For groups $a, b \in \mathcal{A}$, $R_a^{\text{clip}} > R_b^{\text{clip}}$ whenever:

$$\|\mathbf{g}_{D_a}\| \left(p_a - \frac{p_a^2}{2} \right) \geq \frac{5}{2}C + \|\mathbf{g}_{D_b}\| \left(1 + p_b + \frac{p_b^2}{2} \right). \quad (5)$$

Why clipping causes unfairness in DP-SGD

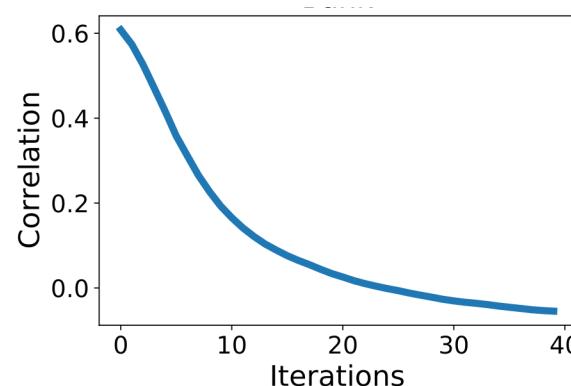
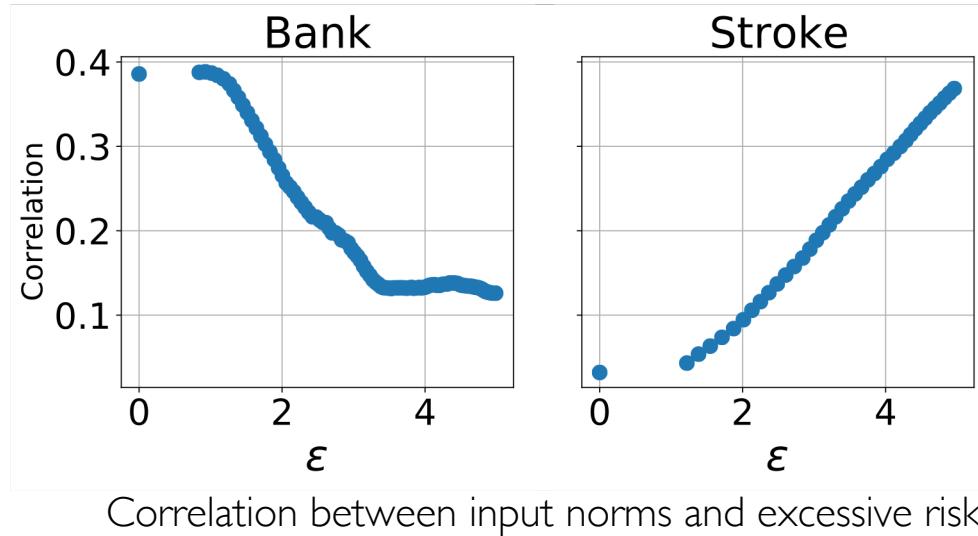


Figure 4: Correlation between inputs and gradients norms.

Crucial proxy to unfairness (due to clipping)

Theorem (informal) Gradients and **Input norms** greatly affect the excessive risk (unfairness) of the individuals and groups

Individuals with **larger features** (e.g., salary, age, height) may have larger disparate impacts in classification!

Why noise addition causes unfairness in DP-SGD?

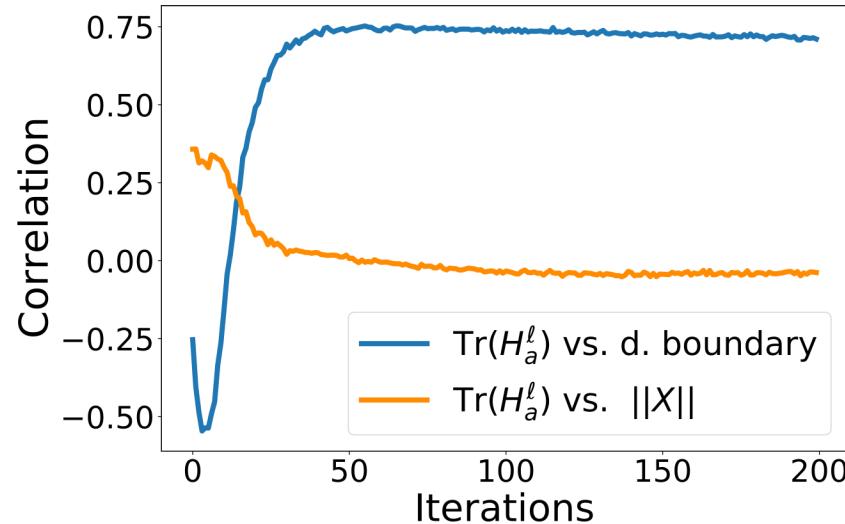


Figure 5: Correlation between trace of Hessian with closeness to boundary (dark color) and input norm (light color).

Crucial proxy to unfairness (due to clipping)

Closeness to the decision boundary!

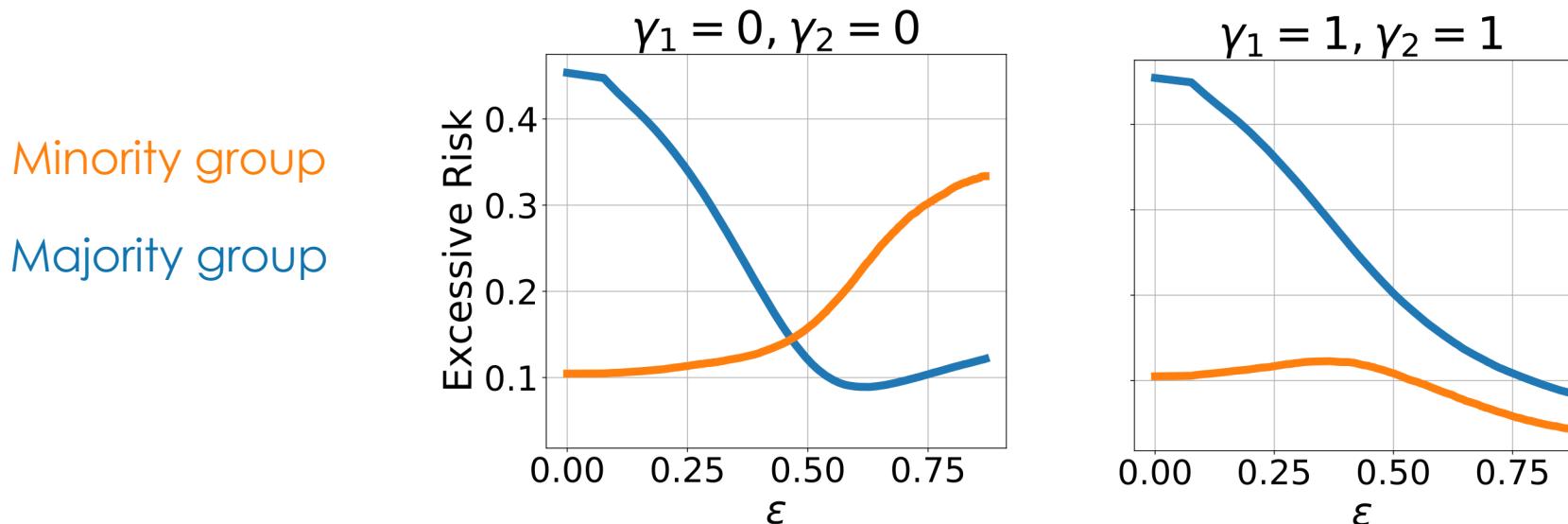
Theorem (informal). Individuals whose outputs are close to the decision boundary will have higher Hessian Traces (high local curvatures of the loss).

Intuitively, the model decisions for samples which are close to the decision boundary are less robust to the presence of noise w.r.t. samples which are farther away from the boundary

Possible mitigating solution

- I. Modify Training so to equalize the factors affecting the excessive risk due to noise and that due to clipping:

$$\min_{\theta} \mathcal{L}(\theta; D) + \sum_{a \in \mathcal{A}} \left(\gamma_1 |\langle g_{D_a} - g_D, g_D - \bar{g}_D \rangle| + \gamma_2 |\text{Tr}(H_\ell^a) - \text{Tr}(H_\ell)| \right),$$



Part III

What are the open scientific challenges?

(Some) Open Questions

- Devise mitigation strategies for data-release mechanisms which account for both fairness and utility.
- Derive more general tools for analyzing fairness for larger classes of decision problems.
- Derive better compositional rules for fairness analysis.
- Analyze the fairness impact of more complex post-processing steps.
- Study the link between Differential Privacy and Robustness under a Fairness Lens.

What are the social and equity risks?

Opinion

Changes to the Census Could Make Small Towns Disappear

By Gus Wezerek and David Van Riper

FEB. 6, 2020







WHY AM I ALWAYS BEING RESEARCHED?

A GUIDEBOOK FOR COMMUNITY ORGANIZATIONS, RESEARCHERS,
AND FUNDERS TO HELP US GET FROM INSUFFICIENT
UNDERSTANDING TO MORE AUTHENTIC TRUTH

Contact Us

nandofioretto@gmail.com



cbowen@urban.org

nandofioretto.com



clairemckaybowen.com

@nandofioretto



@ClaireMKBowen

References

-  Cynthia Dwork, Frank McSherry, Kobbi Nissim, Adam D. Smith: Calibrating Noise to Sensitivity in Private Data Analysis. TCC 2006: 265-284
-  Ferdinando Fioretto, Pascal Van Hentenryck, Keyu Zhu: Differential privacy of hierarchical Census data: An optimization approach. Artif. Intell. 296: 103475 (2021)
-  Keyu Zhu, Ferdinando Fioretto, Pascal Van Hentenryck: Post-processing of Differentially Private Data: A Fairness Perspective. CoRR abs/2201.09425 (2022)
-  David Pujol, Ryan McKenna, Satya Kuppam, Michael Hay, Ashwin Machanavajjhala, Gerome Miklau: Fair decision making using privacy-protected data. FAT* 2020: 189-199
-  Cuong Tran, Ferdinando Fioretto, Pascal Van Hentenryck: Differentially Private and Fair Deep Learning: A Lagrangian Dual Approach. AAAI 2021: 9932-9939
-  Keyu Zhu, Pascal Van Hentenryck, Ferdinando Fioretto: Bias and Variance of Post-processing in Differential Privacy. AAAI 2021: 11177-11184
-  Cuong Tran, Ferdinando Fioretto, Pascal Van Hentenryck, Zhiyan Yao: Decision Making with Differential Privacy under a Fairness Lens. IJCAI 2021: 560-566
-  Cuong Tran, My H. Dinh, Ferdinando Fioretto: Differentially Private Deep Learning under the Fairness Lens. NeurIPS 2021
- Cuong Tran, My H. Dinh, Kyle Beiter, Ferdinando Fioretto: A Fairness Analysis on Private Aggregation of Teacher Ensembles. CoRR abs/2109.08630 (2021)