

# Private Table Statistics using Synthetic Microdata Generation

Terrance D. Savitsky <sup>1</sup>   Matthew R. Williams <sup>2</sup>  
Monika (Jingchen) Hu <sup>3</sup>

<sup>1</sup> U.S. Bureau of Labor Statistics (Office of Survey Methods Research)

<sup>2</sup>RTI International (Division for Statistical and Data Sciences)

<sup>3</sup>Vassar College (Mathematics and Statistics Department)

FCSM 2023 Conference  
October 25, 2023

# Overview

- ▶ **Goal: Generate private data** for tabular release
  - ▶ e.g. Tables of counts and salaries
  - ▶ point estimates and SE estimates
- ▶ **Approach: Synthesize data** with privacy guarantee
  - ▶ Model both outcome  $y$  and survey weights  $w$
  - ▶ Two ways:
    - ▶ 1. Model under observed sample distribution
    - ▶ 2. Model under population distribution
- ▶ **Results:**
  - ▶ Compare synthesizers with additive noise mechanism

# Outline

Differential privacy

Two synthesizers

Laplace Mechanism

Survey of Doctoral Recipients Application

Simulation Study

Concluding remarks

# Differential privacy

- ▶  $D \in \mathbb{R}^{n \times k}$  be a database in input space  $\mathcal{D}$
- ▶ Mechanism  $\mathcal{M}() : \mathbb{R}^{n \times k} \rightarrow \mathcal{O}$ .
- ▶  $\mathcal{M}$  is  $\epsilon$ -differentially private if

$$\frac{\Pr[\mathcal{M}(D) \in O]}{\Pr[\mathcal{M}(D') \in O]} \leq \exp(\epsilon),$$

- ▶ Probability  $\mathcal{M}(D')$  assigns to  $O$  changes by max of  $\exp(\epsilon)$  after deleting 1 row
- ▶ For all  $D, D' \in \mathcal{D}$  that differ by 1 row.

# $\mathcal{M}$ = Additive Noise

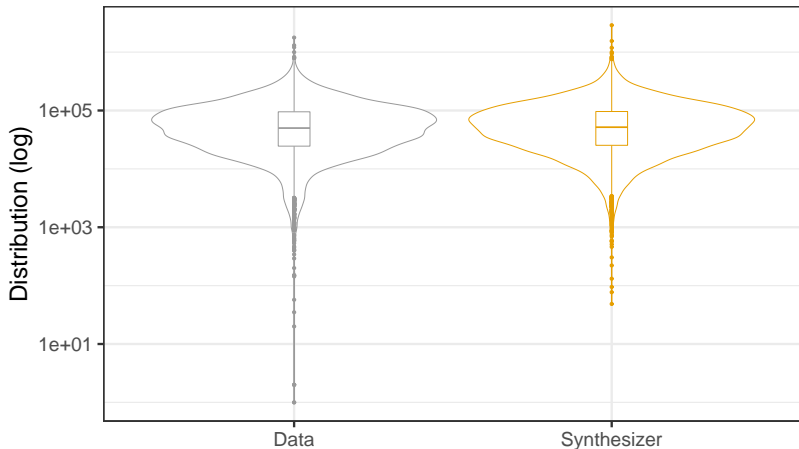
- ▶ An output statistic  $f(D)$ ; e.g., total employment
- ▶ Global sensitivity  
$$\Delta_G = \sup_{D, D' \in \mathcal{D}: \delta(D, D')=1} | f(D) - f(D') |$$
- ▶ Laplace Mechanism for additive noise, scaled to be proportional to  $\Delta_G/\epsilon$  with  $\epsilon$ -DP guarantee
- ▶ Survey weights dramatically increase  $\Delta_G$ .
- ▶ Adding noise disrupts tabular constraints

# $\mathcal{M}$ = Pseudo posterior distribution

$$\xi^{\alpha(\mathbf{y})}(\theta | \mathbf{y}) \propto \prod_{i=1}^n p(y_i | \theta)^{\alpha_i} \times \xi(\theta)$$

- ▶ Down weight each likelihood by  $\alpha_i \in [0, 1]$
- ▶  $\alpha_i$  **lower** when disclosure risk **higher**
- ▶ **Sensitivity** of  $\xi^{\alpha(\mathbf{y})}(\theta | \mathbf{y}) \rightarrow f_{\theta}^{\alpha}(\mathbf{y}) = \log \prod_{i=1}^n p(y_i | \theta)^{\alpha_i}$ .
- ▶  $\Delta_{\alpha} = \sup_{\mathbf{y}, \mathbf{y}' \in \mathcal{Y}^n: \delta(\mathbf{y}, \mathbf{y}')=1} \sup_{\theta \in \Theta} |\alpha(\mathbf{y}) \times f_{\theta}(\mathbf{y}) - \alpha(\mathbf{y}') \times f_{\theta}(\mathbf{y}')|$
- ▶ Posterior draw with  $\epsilon_{\mathbf{y}} = 2\Delta_{\alpha}$  produces one **synthetic**  $\mathbf{y}^*$
- ▶  $\mathbf{y}^*$  produces survey tables with same privacy guarantee

# Synthesizers Encode Privacy by Smoothing



# Asymptotic DP for $\mathcal{M} = \text{Pseudo Posterior}$

To justify a **global** DP result (bounding **all** data sets) compared to a **local** DP result (bound **observed** data set):

**Asymptotic** “Discovery” of  $\Delta_\alpha$  at large sample sizes ( $n$ )

- ▶ Space of **plausible** values  $\Theta$  collapses to a **point**  $\theta^*$ , so don't need to look at  $\sup_{\theta \in \Theta}$ .
- ▶ **Variation** across local  $\Delta_{\alpha, \mathbf{x}}$  **collapses** onto  $\Delta_\alpha$ .
- ▶ **Achieves**  $(\epsilon, \delta)$ -pDP, where  $\delta > 0$  is **probability**  $\exists \mathbf{x} \in \mathcal{X}^n$  **exceeding** the  $\epsilon$  bound.
  - ▶  $\delta \rightarrow 0$  at  $\mathcal{O}(n^{-1/2})$ .
- ▶ **Requires** increasing sparsity in downweighted record contributions, which aligns with focus on isolated records as risky.



# Data from Survey sampling procedure

- ▶ Sample  $S$  of size  $n$  taken from population  $U$  of size  $N$
- ▶ Each individual in  $U$  assigned selection probability  
 $P(\omega_i = 1 \mid \mathcal{A}) = \pi_i$
- ▶ Estimate area statistics with survey weights  $w_i = 1/\pi_i$ , to reduce bias
- ▶ Survey weights,  $(w_i)$ , designed to **correct** bias
- ▶ Incorporating privacy designed to **induce** bias

# Outline

Differential privacy

Two synthesizers

Laplace Mechanism

Survey of Doctoral Recipients Application

Simulation Study

Concluding remarks

## Two Synthesizing Models

- ▶ Synthesis of a local survey database  $(\mathbf{y}_n, \mathbf{w}_n | \mathbf{X}_n, \boldsymbol{\alpha}_n)$ :
- ▶ A Fully Bayes model for **observed sample (FBS)** models  $(\mathbf{y}_n, \mathbf{w}_n | \mathbf{X}_n, \boldsymbol{\alpha}_n)$  under a multinormal pseudo likelihood.
- ▶ A Fully Bayes model for the **population (FBP)** that forms the exact likelihood for  $(\mathbf{y}_n | \mathbf{X}_n, \boldsymbol{\alpha}_n)$ ,  $(\mathbf{w}_n | \mathbf{y}_n, \mathbf{X}_n, \boldsymbol{\alpha}_n)$  in the observed sample.

$$(y_i, w_i | x_i, \alpha_i, \omega_i = 1) = \frac{Pr(\omega_i = 1 | y_i, x_i, w_i) \times (y_i, w_i | x_i, \alpha_i)}{Pr(\omega_i = 1 | x_i, w_i)}$$

- ▶ **FBP** produces synthesized  $\mathbf{y}_n^*$  without sampling bias, so discard weights to build tabular statistics
- ▶ **FBS** requires use of both  $(\mathbf{y}_n^*, \mathbf{w}_n^*)$ .

# Estimation Algorithm

1. Estimate unweighted  $\theta$  with model,

$$\xi(\theta | \mathbf{y}, \mathbf{w}) \propto [\prod_{i=1}^n \pi(y_i, w_i | \theta)] \times \pi(\theta)$$

2. Compute weights,

$$\alpha_i = m(\sup_{\theta \in \Theta} |f_{\theta}(y_i, w_i)|) \propto 1 / \sup_{\theta \in \Theta} |f_{\theta}(y_i, w_i)|$$

3. Re-estimate  $\theta$  using weights,  $\alpha_i$  in

$$\xi^{\alpha}(\theta | \mathbf{y}, \mathbf{w}, \gamma) \propto [\prod_{i=1}^n \pi(y_i, w_i | \theta)^{\alpha_i}] \pi(\theta | \gamma)$$

4. Compute log-likelihood bound,

$$\sup_{y_i, w_i \in \mathcal{D}^n} \sup_{\theta \in \Theta} |\alpha(y_i, w_i) f_{\theta}(y_i, w_i)| \leq \Delta_{\alpha}$$

5. Gives us privacy guarantee,  $\epsilon \leq 2\Delta_{\alpha}$

6. Generate synthetic data,  $(\mathbf{y}^*, \mathbf{w}^*) \sim \pi_{\alpha}(\mathbf{y}^*, \mathbf{w}^* | \mathbf{y}, \mathbf{w})$

# Outline

Differential privacy

Two synthesizers

**Laplace Mechanism**

Survey of Doctoral Recipients Application

Simulation Study

Concluding remarks

# Computing Sensitivity

- ▶ The local sensitivity  $\Delta_{f,g}^c$  for count of field  $f$  and gender  $g$  (cell count):

$$\Delta_{f,g}^c = \max_{i \in \mathcal{S}_{f,g}} w_i - \min_{i \in \mathcal{S}_{f,g}} w_i$$

- ▶ The local sensitivity  $\Delta_{f,g}^a$  for average salary of field  $f$  and gender  $g$  (cell average):

$$\Delta_{f,g}^a = \frac{\max_{i \in \mathcal{S}_{f,g}} w_i y_i - \min_{i \in \mathcal{S}_{f,g}} w_i y_i}{\sum_{i \in \mathcal{S}_{f,g}} w_i - (\max_{i \in \mathcal{S}_{f,g}} w_i - \min_{i \in \mathcal{S}_{f,g}} w_i)}$$

- ▶  $\Delta_*^c = \max_{f,g} \Delta_{f,g}^c$  and  $\Delta_*^a = \max_{f,g} \Delta_{f,g}^a$
- ▶ Generate noise  $\text{Laplace}(0, \Delta_{f,g^*}^{c,a} / \epsilon)$  added to cell count and average salary

# Outline

Differential privacy

Two synthesizers

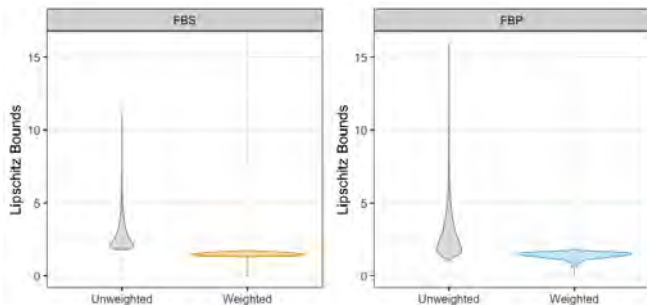
Laplace Mechanism

Survey of Doctoral Recipients Application

Simulation Study

Concluding remarks

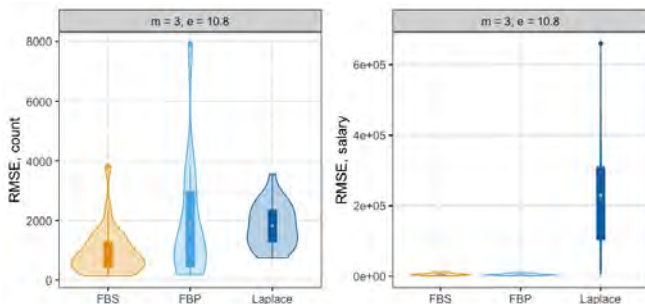
# SDR application (10,355 obs): model fits



**Figure:** Distributions of record-level Lipschitz bounds of the non-private unweighted and the private weighted of FBS (left) and FBP (right) in the SDR application.



# SDR application (10,355 obs): utility evaluation



**Figure:** RMSE values of **counts** (left) and **average salary values** (right) of the three methods, FBS, FBP, and Laplace, applied to the SDR sample. Each violin plot represents a distribution of RMSE values over 27 cells. Results are based on  $m = 3$  synthetic datasets by FBS and FBP, achieving  $\epsilon_{y_n} = 10.8$  for all three methods.

# Outline

Differential privacy

Two synthesizers

Laplace Mechanism

Survey of Doctoral Recipients Application

**Simulation Study**

Concluding remarks

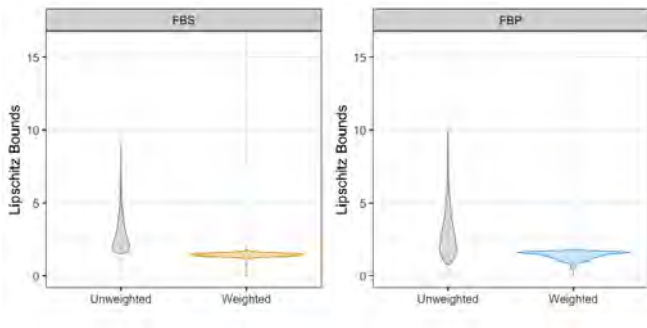
# Simulation studies: simulation design

- ▶ Based on the 2017 SDR public use file
- ▶ Population  $N = 100,000$  units of: salary ( $y_i$ ), field of expertise and gender ( $x_i$ )
- ▶ Salary  $y_i \mid x_i \sim \text{Lognormal}(\mu_i, 0.4)$  where  $\mu_i$  is group-specific mean from the public use file
- ▶ Additive noise:  $\text{noise}_i \sim \text{Lognormal}(0, 0.4)$
- ▶ Survey weights:

$$\begin{aligned}\log(\pi_i) &= \log(y_i) + \text{noise}_i \\ w_i &= 1/\pi_i\end{aligned}$$

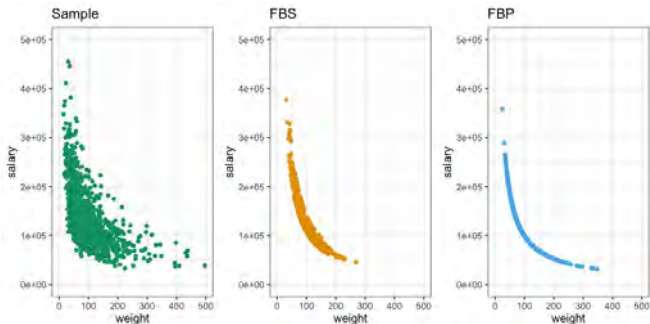
- ▶ Take a stratified PPS sample of  $n = 1000$  using fields as strata:  $(y_n, X_n, w_n)$

# Sensitivity before/after weighting by $\alpha$



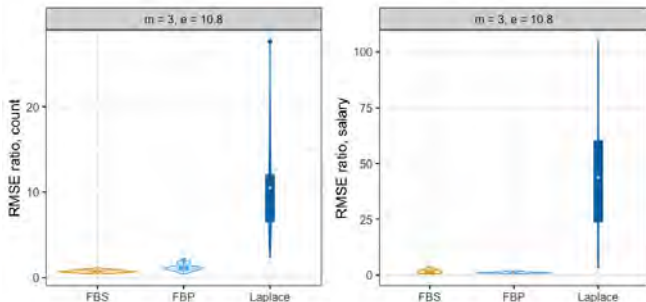
**Figure:** Distributions of record-level sensitivity bounds of the non-private unweighted and the private weighted of FBS (left) and FBP (right) in the simulation.

# Smoothed weights $\mathbb{E}(w|y)$ Improves Efficiency



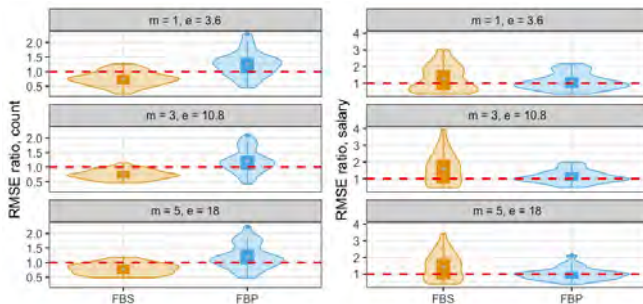
**Figure:** Comparison of the salary and weight bivariate distributions of confidential salary and weights in the sample (green and left), synthetic salary and smoothed weights from FBS (yellow and middle), and synthetic salary and smoothed weights from FBP (blue and right).

# Utility of $\mathcal{M} = \text{Synthesizers}$ versus $\mathcal{M} = \text{Laplace}$



**Figure:** RMSE ratios of **counts** (left) and **average salary values** (right) of the three methods, FBS, FBP, and Laplace, applied to the selected sample. Each violin plot represents a distribution of RMSE ratios over 27 cells. Results are based on  $m = 3$  synthetic datasets by FBS and FBP, achieving  $\epsilon_{y_n} = 10.8$  for all three methods.

# Adding Synthetic Data Replicates, $m$ , Improves Utility



**Figure:** RMSE ratios of **counts** (left) and **average salary values** (right) of FBS and FBP, applied to the selected sample. A red dashed line at RMSE ratio = 1 is included for reference. Each violin plot represents a distribution of RMSE ratios over 27 cells. Results are based on  $m = \{1, 3, 5\}$  synthetic datasets by FBS and FBP, achieving  $\epsilon_{y_n} = \{3.6, 10.8, 18\}$  for both methods.

# Outline

Differential privacy

Two synthesizers

Laplace Mechanism

Survey of Doctoral Recipients Application

Simulation Study

Concluding remarks



# Summary

- ▶ Formal privacy for data collected under an informative sampling design
- ▶ We recommend the FBS: easy to estimate and produces low RMSE
- ▶ The synthetic data is privacy protected and obeys all constraints without any post processing
- ▶ No interactive queries required
- ▶ The synthetic data may be used for other purposes
- ▶ arXiv link to manuscript:  
<https://arxiv.org/abs/2101.06188>

# References

- ▶ Leon-Novelo, L. G. and Savitsky, T. D. (2019). Fully Bayesian estimation under informative sampling, *Electronic Journal of Statistics*, 13, 1608-1645.
- ▶ Savitsky, T. D., Williams, M. R. and Hu, J. (2020). Bayesian pseudo posterior mechanism under differential privacy. arXiv:1909.11796.
- ▶ Rao, J. N. K., Wu, C. F. J. and Yue, K. (1992). Some recent work on resampling methods for complex surveys, *Survey Methodology*, **18**, 209-217.