

# Privacy Amplification for Synthetic data using Range Restriction

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

<sup>1</sup> Binghamton University (Department of Mathematics and Statistics)

<sup>2</sup> RTI International (Center for Official Statistics)

<sup>3</sup> U.S. Bureau of Labor Statistics (Office of Survey Methods Research)

2026 Privacy and Public Policy Conference



# Outline

Synthetic Data and Formal Privacy

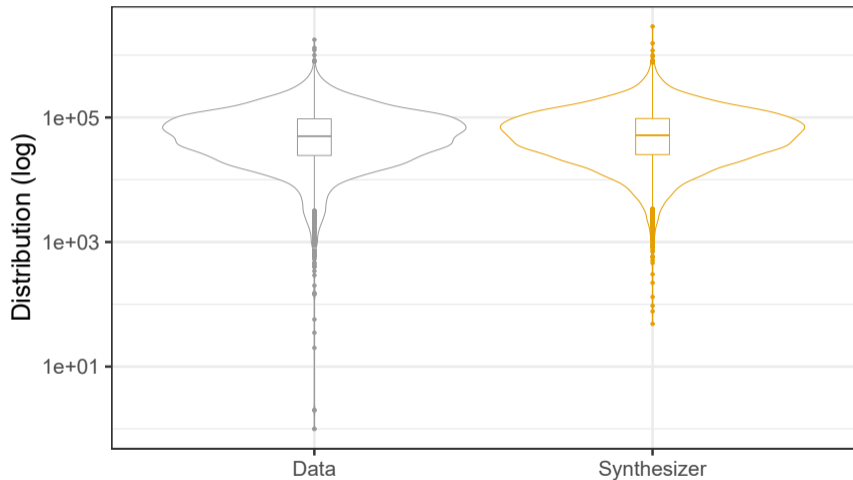
Pseudo Posterior Mechanism

Privacy Amplification by Conditioning on Known Subspaces

# Synthetic data

- ▶ Respondent-level data generation and release
- ▶ Estimate a model on the confidential microdata
- ▶ Simulate synthetic data from posterior predictive distribution
- ▶ Constructing tables is one way to measure synthetic data utility
- ▶ **Unlimited queries**; e.g. survey tables, regression analyses etc.
- ▶ High utility and low risks (**non-formal**)

# Replicate data $y^*$ given (|) observed data, $x$ (Hu, 2019)



# Formal Privacy

- ▶ **Formal Privacy** is a property associated with release mechanisms  $\mathcal{M}$  (e.g. generating a synthetic data set).
  - ▶ has a mathematical **bound** on the **information** contribution of individuals
  - ▶ allows for **privacy accounting** (e.g. summing privacy loss/expense across all uses of the data)
  - ▶ requires that the mechanism has to be **random** (stochastic)
- ▶ The most popular formal privacy definition (family) is **differential privacy (DP)**, where the bounds are over neighboring data sets with and without individual data

# Why Differential Privacy (DP)?

## In General

- ▶ DP links amount of randomization to risk thresholds.
- ▶ DP calls for a rigorous tracking of all data uses and releases.
- ▶ DP does not make explicit assumptions about the knowledge and behavior of the intruder.

## For Our Discussion

- ▶ Additivity of risk across releases based on worst case scenario .
- ▶ Room for improvement, especially if worst case is  $\infty$  risk
- ▶ Room for improvement, especially if some knowledge is considered public.

# Outline

Synthetic Data and Formal Privacy

Pseudo Posterior Mechanism

Privacy Amplification by Conditioning on Known Subspaces

# Differential Privacy under $\mathcal{M} = \xi(\theta | \mathbf{x})$

$$\sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n: \delta(\mathbf{x}, \mathbf{x}')=1} \sup_{B \in \beta_{\Theta}} \frac{\xi(B | \mathbf{x})}{\xi(B | \mathbf{x}')} \leq e^{\epsilon},$$

- ▶  $\epsilon$  bounds the **change** in the **probability measure**  $\xi$ 
  - ▶ from the inclusion of a **single record**  $\delta(\mathbf{x}, \mathbf{x}') = 1$ ,
  - ▶ over **all possible outcomes**,  $B \in \beta_{\Theta}$  – sets in the space of measurable sets of  $\Theta$ .
  - ▶ over **all possible data sets**  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$  of size  $n$ .

# DP under the Posterior Mechanism

Dimitrakakis et al. (2017) establish a link between bounding the log-likelihood  $f_{\theta}(\mathbf{x}) = \log \pi_{\theta}(\mathbf{x})$  and a DP bound  $\epsilon$ .

- ▶ If  $\sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n: \delta(\mathbf{x}, \mathbf{x}')=1} \sup_{\theta \in \Theta} |f_{\theta}(\mathbf{x}) - f_{\theta}(\mathbf{x}')| \leq \Delta$ .
- ▶ Then a single draw of replicate data  $\mathbf{x}^* \sim \pi(\mathbf{x}^* | \mathbf{x})$  has a DP guarantee of  $\epsilon \leq 2\Delta$ .

# Major Limitations of the Posterior Mechanism

The results in Dimitrakakis et al. (2017) are **not immediately** practical.

- ▶ For many data distributions of interest (e.g. normal, exponential, Poisson, geometric),  $\Delta = \infty$ .
  - ▶ Directly **truncating** the support  $\Theta$  and the domain  $\mathcal{X}^n$  may work for simple distributions but is somewhat **ad-hoc** and **scales poorly** with the dimension of both
- ▶ Using a **prior** that induces more **smoothing** may reduce  $\Delta$  – but requires **re-estimation**. This is an **indirect** adjustment.

# Pseudo Posterior Mechanism

- ▶ Savitsky et al. (2022) utilize record-indexed weights,  $\alpha \in (0, 1]^n$
- ▶ To **downweight** likelihood contributions with **high disclosure risk**

$$\xi^\alpha(\theta | \mathbf{x}, \gamma) \propto \left[ \prod_{i=1}^n \pi(x_i | \theta)^{\alpha_i} \right] \pi(\theta | \gamma)$$

- ▶  $\alpha_i \propto 1 / \sup_{\theta \in \Theta} |f_\theta(x_i)|$
- ▶ Allows **surgical** downweighting of high risk records
- ▶  $\alpha_i$  induces an anti-informative prior
- ▶  $\sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n: \delta(\mathbf{x}, \mathbf{x}')=1} \sup_{\theta \in \Theta} |\alpha(\mathbf{x})f_\theta(\mathbf{x}) - \alpha(\mathbf{x}')f_\theta(\mathbf{x}')| \leq \Delta_\alpha$
- ▶  $\Delta_\alpha \leq \Delta$
- ▶ **Ensures**  $\Delta_\alpha < \infty$
- ▶ Expected to better preserve real data distribution for any target privacy budget,  $\epsilon$

# Outline

Synthetic Data and Formal Privacy

Pseudo Posterior Mechanism

Privacy Amplification by Conditioning on Known Subspaces

## (A) Measure Expert Beliefs as modeled Probability

- ▶ Expert believes interested public already knows that the true datum value does not lie in the space outside of the ball/interval
- ▶  $p_{\mathcal{M}}(\mathbf{x}^* | \mathbf{x})$ , constructed from the same model estimated on the first step to construct the risk-based weights,  $\alpha_i$ .

$$\lambda_i = \Pr_{\mathcal{M}, a, b} (x_i^* \notin [a \times x_i, b \times x_i]) = \Pr_{\mathcal{M}, R_i} (x_i^* \notin R_i)$$

- ▶ Each  $\lambda_i$ , then, represents the probability that a realized datum value for unit  $i$  does *not* need to be protected

## (A) Risk-based Weights Adjusted by Expert Beliefs

$$\alpha_i^* = \lambda_i + (1 - \lambda_i) \times \alpha_i$$

- ▶ When a **unit is highly risky** such that  $\alpha_i = 0$  (is set to 0) then  $\alpha_i^* = \lambda_i$ .
- ▶ By contrast, when datum  $i$  presents **no disclosure risk** such that  $\alpha_i = 1$ , then  $\alpha_i^* = 1$ , indicating there is no privacy protection required under this scenario.
- ▶ When  $\lambda_i = 0$  it is supposed that the interested public has **no knowledge** about a subspace of the private variable and we return to the usual case where  $\alpha_i^* = \alpha_i$

## (A) Reduced Sensitivity by Incorporating Probabilities

$$p_{\theta}^{\alpha^*}(\mathbf{x}) = p_{\theta}^{\lambda^c \alpha}(\mathbf{x}) p_{\theta}^{\lambda}(\mathbf{x}) = \prod_{i=1}^n p(x_i | \theta)^{(1-\lambda_i)\alpha_i} \prod_{i=1}^n p(x_i | \theta)^{\lambda_i}.$$

$$\begin{aligned} \Delta_{\alpha, \lambda, \mathbf{x}} &= \max_{\theta \in \{\xi^{\alpha^*}(\mathbf{x}) | \theta | \mathbf{x}\}_m} \max_{i \in 1, \dots, n} |\alpha_i^* \times f_{\theta_m}(x_i) - \lambda_i \times f_{\theta_m}(x_i)| \\ &= \max_{\theta \in \{\xi^{\alpha^*}(\mathbf{x}) | \theta | \mathbf{x}\}_m} \max_{i \in 1, \dots, n} |((1 - \lambda_i)\alpha_i) \times f_{\theta_m}(x_i)| \leq \Delta_{\alpha, \mathbf{x}}, \end{aligned}$$

- ▶ Further reduce the sensitivity by factor of  $(1 - \lambda_i)$ .
- ▶ The  $(1 - \lambda_i)$  term represents the **portion of the data support that is *not* known** by the interested public and, therefore, sensitive.
- ▶ Under our set-up, we account for expert beliefs *indirectly* through knowledge weights  $\lambda$  rather than directly truncating the data support,  $\mathcal{X}^n$ .

## (B) Minimal Use of Expert Knowledge

- ▶ The specification of the range  $R_i = [a \times x_i, b \times x_i]$  gives rise to
  - ▶ the truncation adjustment  $P_\theta(R_i) = \int_{x \in R_i} p_\theta(x) dx$
  - ▶ the **range-truncated likelihood**  $p_{\theta_i}(x_i)/P_\theta(R_i)$
- ▶ Assuming the range  $R_i$  is public, we can use  $R_i$  for 'free' in inference, analogous to using censoring times in survival analysis.

$$\begin{aligned} p^{I^c, \alpha_i}(x_i | \theta, a, b) &= p(x_i | \theta)^{\alpha_i} / (P(b \times x_i | \theta) - P(a \times x_i | \theta)) \\ &= p(x_i | \theta)^{\alpha_i} / p^I(x_i | \theta, a, b) \end{aligned}$$

## (B) Reduced Sensitivity by Incorporating End Points

$$p_{\theta}^{\alpha}(\mathbf{x}) = p_{\theta}^{I^c, \alpha}(\mathbf{x}) p_{\theta}^I(\mathbf{x}) = \prod_{i=1}^n p^{I^c, \alpha_i}(x_i | \theta, a, b) \prod_{i=1}^n p^I(x_i | \theta, a, b).$$

$$\Delta_{\alpha, I, \mathbf{x}} = \max_{\theta \in \{\xi(\theta | \mathbf{x})\}_m} \max_{i \in \{1, \dots, n\}} |\alpha_i f_{\theta_m}(x_i) - \log(P_{\theta_m}(R_i))| \leq \Delta_{\alpha, \mathbf{x}}$$

- ▶ Further reduce the risk-based weighted sensitivity by

$$\log(P_{\theta_m}(R_i)) = \log(P(b \times x_i | \theta_m) - P(a \times x_i | \theta_m))$$

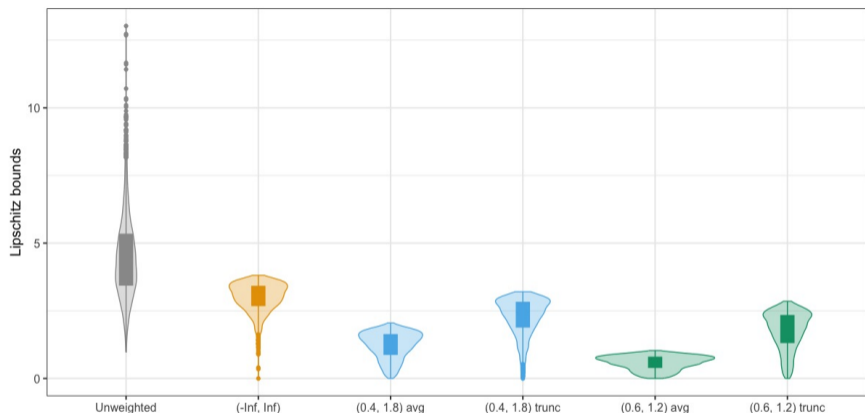
- ▶ Rather than directly **truncating the data support**,  $\mathcal{X}^n$  and greatly increasing the complexity of the synthesizer  $p_{\theta}^{\alpha}(\mathbf{x})$ , we **partition the sensitivity** calculation into contributions coming from within and outside the restricted range.

# Range Restricted Privacy under the Pseudo Posterior

$$(A) \quad \sup_{\mathbf{x} \in \mathcal{X}^n, \mathbf{x}' \in \mathcal{X}^{n-1}: \delta(\mathbf{x}, \mathbf{x}')=1} \sup_{B \in \beta_{\Theta}} \frac{\xi^{\lambda^c \alpha(\mathbf{x})}(B \mid \mathbf{x}, \mathbf{R})}{\xi^{\lambda^c \alpha(\mathbf{x})}(B \mid \mathbf{x}, \mathbf{R})} \leq e^{\epsilon},$$

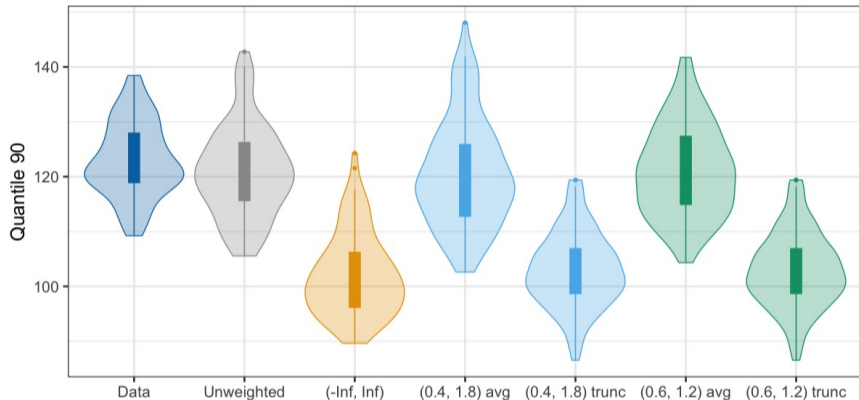
$$(B) \quad \sup_{\mathbf{x} \in \mathcal{X}^n, \mathbf{x}' \in \mathcal{X}^{n-1}: \delta(\mathbf{x}, \mathbf{x}')=1} \sup_{B \in \beta_{\Theta}} \frac{\xi^{I^c \alpha(\mathbf{x})}(B \mid \mathbf{x}, \mathbf{R})}{\xi^{I^c \alpha(\mathbf{x}')}(B \mid \mathbf{x}', \mathbf{R})} \leq e^{\epsilon},$$

# Sensitivities under Intruder Knowledge



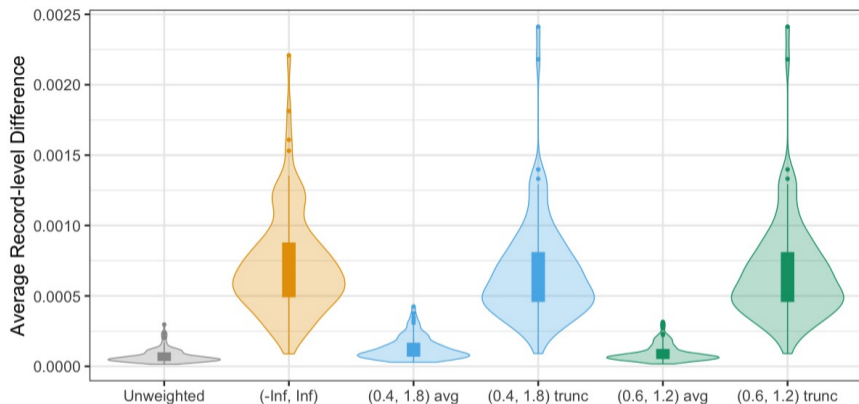
**Figure:** Violin plots of by-record Lipschitz bounds of Unweighted, (-Inf, Inf) (i.e., no bounds as Weighted), (0.4, 1.8) averaged, (0.4, 1.8) truncated, (0.6, 1.2) averaged, and (0.6, 1.2) truncated, over a single sample.

## Preservation of Quantiles under MC study



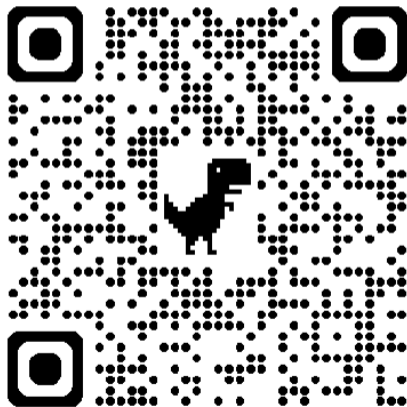
**Figure:** Violin plots of Q90s of Unweighted,  $(-\text{Inf}, \text{Inf})$  (i.e., no bounds as Weighted),  $(0.4, 1.8)$  averaged,  $(0.4, 1.8)$  truncated,  $(0.6, 1.2)$  averaged, and  $(0.6, 1.2)$  truncated, over 100 repeated samples.

# Real-vs-Synthetic Distribution Distance under MC study



**Figure:** Violin plots of average metric of ECDF of Unweighted,  $(-\infty, \infty)$  (i.e., no bounds as Weighted),  $(0.4, 1.8)$  averaged,  $(0.4, 1.8)$  truncated,  $(0.6, 1.2)$  averaged, and  $(0.6, 1.2)$  truncated, over 100 repeated samples.

<https://arxiv.org/abs/2602.04124>



# References I

- Dimitrakakis, C., Nelson, B., Zhang, Z., Mitrokotsa, A. and Rubinstein, B. I. P. (2017), 'Differential privacy for bayesian inference through posterior sampling', *Journal of Machine Learning Research* **18**(1), 343–381.
- Hu, J. (2019), 'Bayesian estimation of attribute and identification disclosure risks in synthetic data', *Transactions on Data Privacy* **12**, 61–89.
- Savitsky, T. D., Williams, M. R. and Hu, J. (2022), 'Bayesian pseudo posterior mechanism under asymptotic differential privacy', *Journal of Machine Learning Research* **23**, 1–37.

# Extra Slides

# Local vs. Global Privacy Guarantee

When we **implement**

- ▶ We have to **estimate**  $\Delta_{\alpha, \mathbf{x}}$  based on a **single** data set,  $\mathbf{x}$ , to estimate the DP guarantee  $\epsilon$ . (local DP result)
- ▶ We have to **approximate** the  $\sup_{\theta \in \Theta} f_{\theta}(x_i)$  as  $\max_{\theta_j, j \in 1, \dots, J} f_{\theta_j}(x_i)$ .
- ▶ **Overestimate**  $\alpha \rightarrow$  **Underestimate**  $\epsilon$

# Asymptotic Global Guarantee

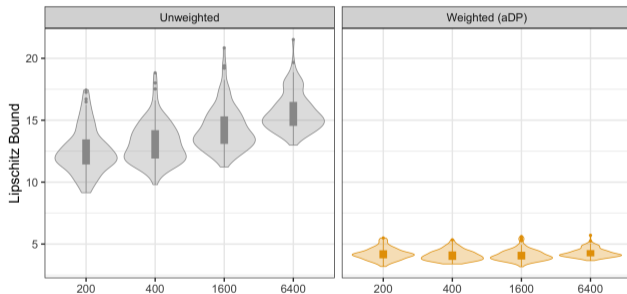
To justify a **global** DP result (bounding **all** data sets) compared to a **local** DP result (bound **observed** data set):

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

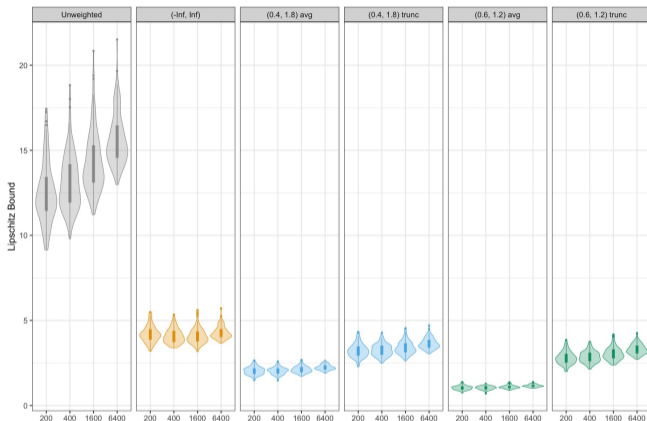
- ▶ Space of **plausible** values  $\Theta$  collapses to a **point**  $\theta^*$ , so don't need to look at  $\sup_{\theta \in \Theta}$ .
- ▶ **Variation** across local  $\Delta_{\alpha, \mathbf{x}}$  **collapses** onto  $\Delta_\alpha$ .
- ▶ **Achieves**  $(\epsilon, \delta)$ -pDP, where  $\delta > 0$  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.

# Contraction of Lipschitz, ( $\Delta_{\alpha, x_r}$ )

- ▶ Generate  $x_{ri} \stackrel{\text{ind}}{\sim}$  lognormal for  $r = 100$ .



# Asymptotic Contraction of Local Privacy Guarantee



**Figure:** Violin plots of Lipschitz bounds of Unweighted,  $(-\text{Inf}, \text{Inf})$  (i.e., no bounds as Weighted),  $(0.4, 1.8)$  averaged,  $(0.4, 1.8)$  truncated,  $(0.6, 1.2)$  averaged, and  $(0.6, 1.2)$  truncated, over 100 repeated samples and four different sample sizes  $n = \{200, 400, 1600, 6400\}$ .