

Bayesian Pseudo Posterior Mechanism under Asymptotic Differential Privacy

Terrance D. Savitsky ¹ Matthew R. Williams ²
Monika (Jingchen) Hu ³

¹ U.S. Bureau of Labor Statistics (Office of Survey Methods Research)

²NSF (National Center for Science and Engineering Statistics)

³Vassar College (Mathematics and Statistics Department)

JSM - 2024 August 7, 2024

Outline

Differential Privacy

Synthetic Data

Pseudo Posterior Mechanism

Construction of Risk-based Weights

Achieving Global DP Guarantee

Simulation Study to Illustrate Asymptotic Contraction of Lipschitz

Applications to the Consumer Expenditure Surveys

Re-weighting Strategy

Posterior Predictive Distribution

- ▶ Distribution of **replicate** data \mathbf{x}^* given (|) **observed** data, \mathbf{x}
- ▶ Used for **imputation** or to generate **synthetic** (or “fake”) data
- ▶ Model-based synthetic data **smooths** real data distribution
- ▶ Smoothing encodes privacy while preserving essential properties of observed data distribution
- ▶ User allowed unlimited ‘interrogation’ of synthetic data

$$\pi(\mathbf{x}^*|\mathbf{x}) = \int \prod_{i=1}^n \pi(x_i^*|\boldsymbol{\theta}) \times \xi(\boldsymbol{\theta}|\mathbf{x}) d\boldsymbol{\theta} \approx \frac{1}{S} \sum_{s=1}^S \prod_{i=1}^n \pi(x_i^*|\boldsymbol{\theta}_s)$$

- ▶ Where $s = 1, \dots, S$ indexes the number of MCMC draws of $\boldsymbol{\theta}$
- ▶ We call this model (distribution) for generating data a **synthesizer**

\mathcal{M} = Posterior distribution

$$\xi(\theta | \mathbf{x}) \propto \prod_{i=1}^n \pi(x_i | \theta) \times \xi(\theta)$$

- ▶ Sensitivity of $\xi(\theta | \mathbf{x})$ based on $f_{\theta}(\mathbf{x}) = \log \prod_{i=1}^n \pi(x_i | \theta)$.
- ▶ $\Delta = \sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n: \delta(\mathbf{x}, \mathbf{x}')=1} \sup_{\theta \in \Theta} |f_{\theta}(\mathbf{x}) - f_{\theta}(\mathbf{x}')|$
- ▶ Generate a **single** dataset under $\epsilon = 2\Delta$ given a draw of θ
- ▶ Each posterior draw with $\epsilon = 2\Delta$ produces one synthetic \mathbf{x}^*

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

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

- ▶ ϵ bounds the **change** in the **probability measure** ξ
 - ▶ 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 Θ .
 - ▶ over **all possible data sets** $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$ of size n .

Outline

Differential Privacy

Synthetic Data

Pseudo Posterior Mechanism

Construction of Risk-based Weights

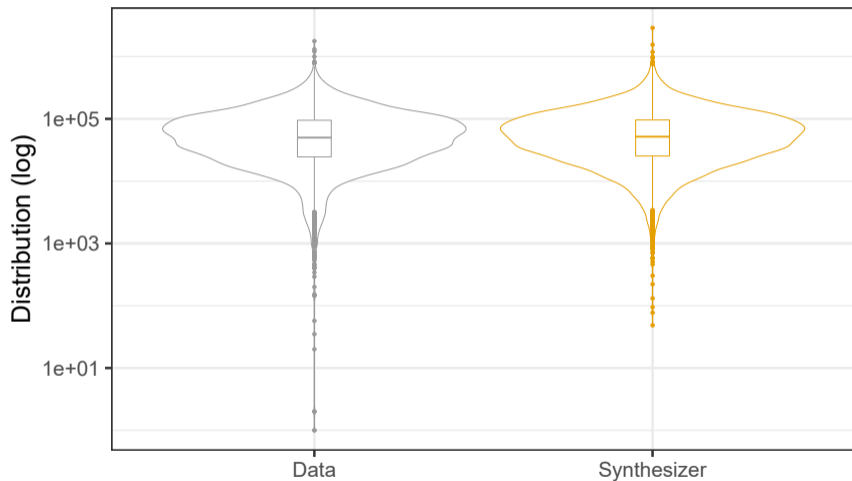
Achieving Global DP Guarantee

Simulation Study to Illustrate Asymptotic Contraction of Lipschitz

Applications to the Consumer Expenditure Surveys

Re-weighting Strategy

Example of a Flexible Synthesizer (Hu and Savitsky, 2019)



Outline

Differential Privacy

Synthetic Data

Pseudo Posterior Mechanism

Construction of Risk-based Weights

Achieving Global DP Guarantee

Simulation Study to Illustrate Asymptotic Contraction of Lipschitz

Applications to the Consumer Expenditure Surveys

Re-weighting Strategy

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** ϵ .

- ▶ 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 Θ 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 Δ – but requires **re-estimation**. This is an **indirect** adjustment.

Pseudo Posterior Mechanism

- ▶ Savitsky et al. (2019) 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
- ▶ α_i induces an anti-informative prior
- ▶ **Ensures** $\Delta_\alpha < \infty$
- ▶ Expected to better preserve real data distribution for any target privacy budget, ϵ

Formulating Risk-Based Weights

How might we choose $\alpha \in (0, 1]^n$?

$$\alpha_i = m \left(\sup_{\theta \in \Theta} |f_{\theta}(x_i)| \right), \quad (1)$$

where $f_{\theta}(x_i)$ is **unweighted** log-likelihood **contribution** from x_i , and $m(z)$ is a **decreasing** function $m : R^+ \rightarrow [0, 1]$. ($m(0) = 1$, $m(\infty) = 0$)

▶ **Riskier** values **downweighted** more

▶ α_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** $\Delta_{\alpha} < \infty$

▶ Produces DP budget of $\epsilon \leq 2\Delta_{\alpha}$.

▶ Control ϵ **indirectly** through setting of $(\alpha_1, \dots, \alpha_n)$

Estimation Algorithm

1. Estimate unweighted θ with model,
 $\xi(\theta|\mathbf{x}) \propto [\prod_{i=1}^n \pi(x_i|\theta)] \times \pi(\theta)$
2. Compute weights, $\alpha_i = m(\sup_{\theta \in \Theta} |f_{\theta}(x_i)|) \propto 1/\sup_{\theta \in \Theta} |f_{\theta}(x_i)|$
3. Re-estimate θ using weights, α_i in
 $\xi^{\alpha}(\theta|\mathbf{x}, \gamma) \propto [\prod_{i=1}^n \pi(x_i|\theta)^{\alpha_i}] \pi(\theta|\gamma)$
4. Compute log-likelihood bound,
 $\sup_{\theta \in \Theta} |\alpha_i \times f_{\theta}(x_i)| \leq \Delta_{\alpha, x_i}$
 $\max_{x_i \in \mathbf{x}^n} \Delta_{\alpha, x_i} \leq \Delta_{\alpha, \mathbf{x}}$
5. Gives us privacy guarantee, $\epsilon \leq 2\Delta_{\alpha, \mathbf{x}}$
6. Generate synthetic data $\mathbf{x}^* \sim \pi_{\alpha}(\mathbf{x}^*|\mathbf{x})$

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 ϵ . (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** ϵ

Ensuring Global DP

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

- ▶ **Asymptotic** Discovery - For large sample sizes (n)
 - ▶ Space of **plausible** values Θ collapses to a **point** θ^* , so don't need to look at $\sup_{\theta \in \Theta}$.
 - ▶ **Variation** across local $\Delta_{\alpha, x}$ **collapses** onto Δ_{α} .
 - ▶ **Achieves** (ϵ, δ) - pDP, where $\delta > 0$ is **probability** $\exists \mathbf{x} \in \mathcal{X}^n$ **exceeding** the ϵ 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.
- ▶ **Truncating** the weights, $(\alpha_1, \dots, \alpha_n)$
 - ▶ **Set** $\alpha_i^* = 0$ if $\alpha_i \times f_{\theta}(x_i) > \Delta_{\alpha}$
 - ▶ Target, global threshold, Δ_{α} , regulates downweighting and speed of contraction.

Outline

Differential Privacy

Synthetic Data

Pseudo Posterior Mechanism

Construction of Risk-based Weights

Achieving Global DP Guarantee

Simulation Study to Illustrate Asymptotic Contraction of Lipschitz

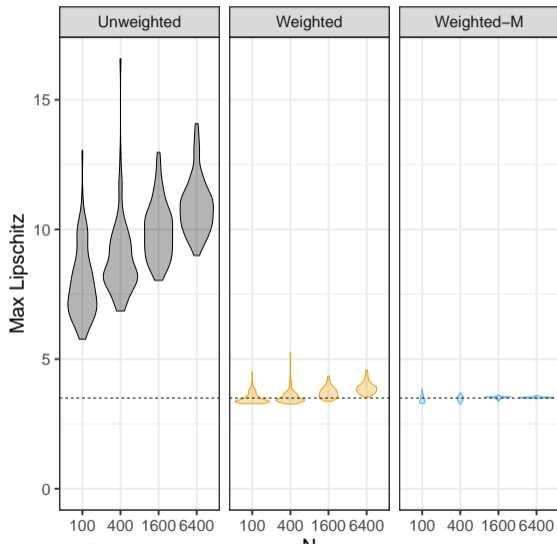
Applications to the Consumer Expenditure Surveys

Re-weighting Strategy

Monte Carlo Simulation at Increasing n

1. For **sample size**, $n \in \{100, 400, 1600, 6400\}$, **repeat** the following Monte Carlo procedure to generate a distribution of local Lipschitz bounds:
2. For $r = 1, \dots, 100$:
 - 2.1 **Generate** $\mathbf{x}_r \sim \text{Pois}(\mu = 100)$, each of size n .
 - 2.2 **Compute** the **local Sensitivity bound**, $\Delta_{\alpha, \mathbf{x}_r}$, for the unweighted, α -weighted, and **M -truncation-weighted pseudo posterior mechanisms**.
 - 2.3 Construct the distribution of $\Delta_{\alpha, \mathbf{x}_r}$ and note the maximum and difference between the maximum and minimum values of the distribution.
3. **Assess contraction** of the $\max_r \Delta_{\alpha, \mathbf{x}_r}$ to a **single (global) value** and whether the minimum and maximum values collapse together.

Contraction of Weighted- M Lipschitz, (Δ_{α, x_r})



Weighted- M Utility Performance

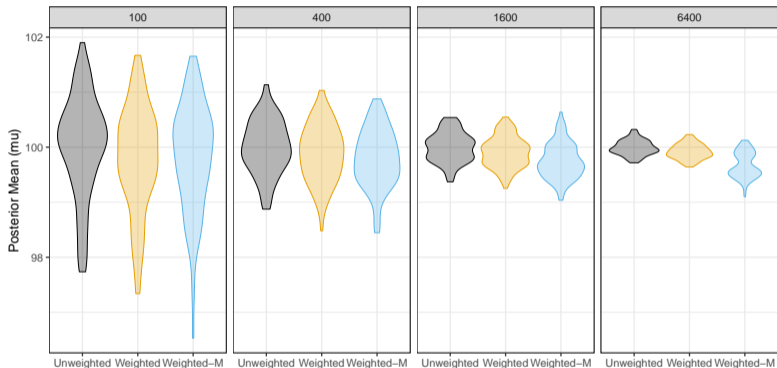


Figure: Distributions of the average of mean parameter μ for each of sample size (100, 400, 1600, 6400) from $R = 100$ realizations.

Outline

Differential Privacy

Synthetic Data

Pseudo Posterior Mechanism

Construction of Risk-based Weights

Achieving Global DP Guarantee

Simulation Study to Illustrate Asymptotic Contraction of Lipschitz

Applications to the Consumer Expenditure Surveys

Re-weighting Strategy

Application: The Consumer Expenditure Surveys

- ▶ Provides **income** and **expenditure** patterns among consumer units (CUs)
- ▶ CE **releases** public-use **microdata** (PUMS)
- ▶ Focus on **sensitive**, skewed **income** variable – with non-sensitive variables as predictors
- ▶ From 1st QTR of 2017 survey of **n = 6208** CUs

Variables in the CE sample

Gender

Age

Education Level

Region

Urban

Marital Status

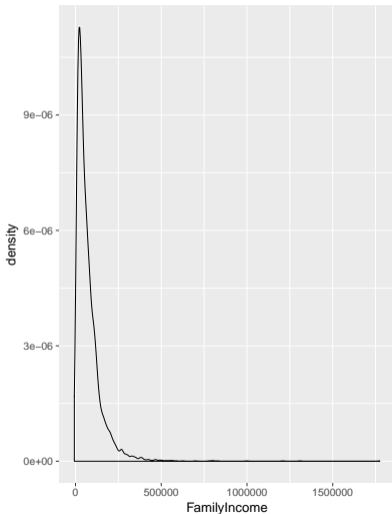
Urban Type

CBSA

Family Size

Earner

****Family Income****



Comparing Synthesizer Bounds by Record ($\epsilon = 2\Delta$)

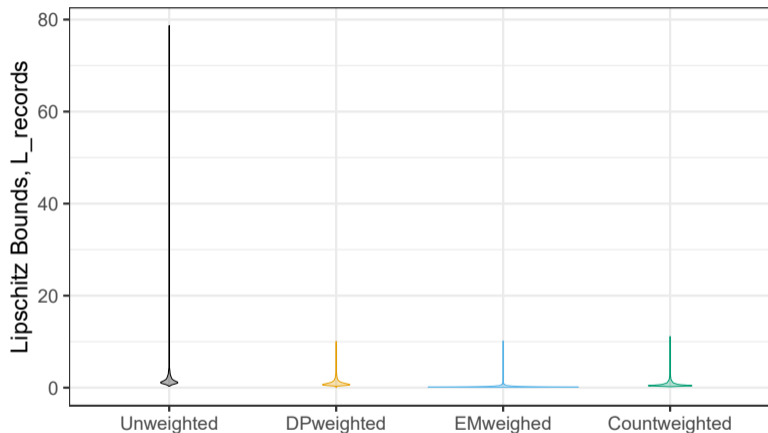
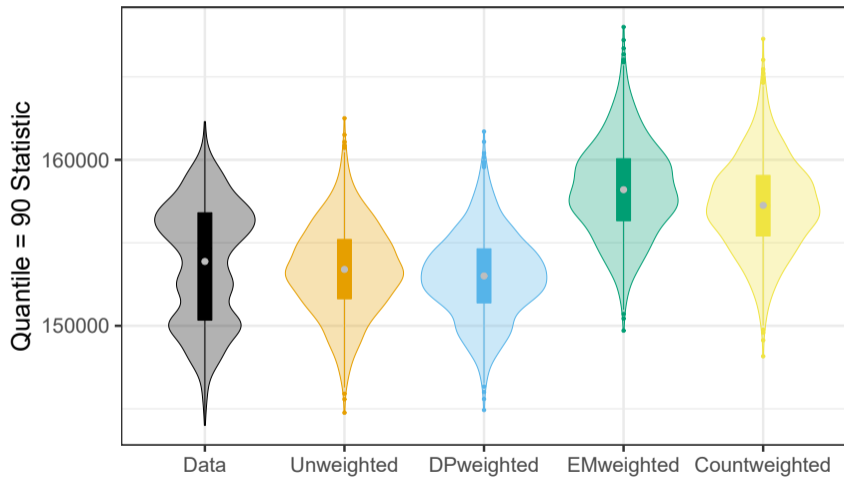


Figure: $\Delta_{Unweighted} = 78.7$, $\Delta_{DP} = 10.1$, $\Delta_{EM} = 10.2$, $\Delta_{Count} = 11.17$.

Comparing Synthesizer Utility (Q90) of Synthetic Data



Outline

Differential Privacy

Synthetic Data

Pseudo Posterior Mechanism

Construction of Risk-based Weights

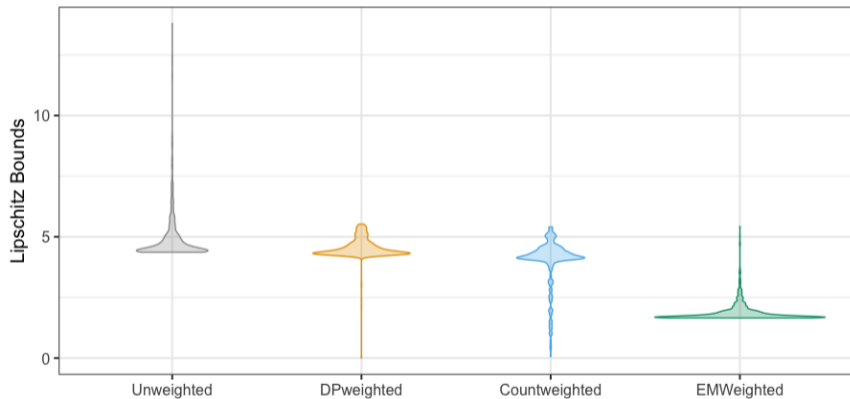
Achieving Global DP Guarantee

Simulation Study to Illustrate Asymptotic Contraction of Lipschitz

Applications to the Consumer Expenditure Surveys

Re-weighting Strategy

Highly Skewed Count Data



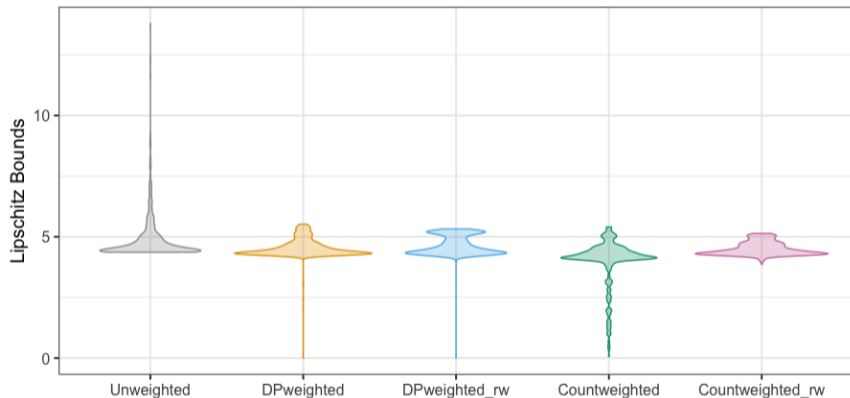
Re-weighting to compress Lipschitz Distribution

- ▶ Formal privacy guarantee is a **supremum** (not an average)
- ▶ There is no **credit** for records whose $\Delta_{\alpha, x_i} < \Delta_{\alpha, x}$
- ▶ Up-weight records with $\Delta_{\alpha, x_i} < \Delta_{\alpha, x}$

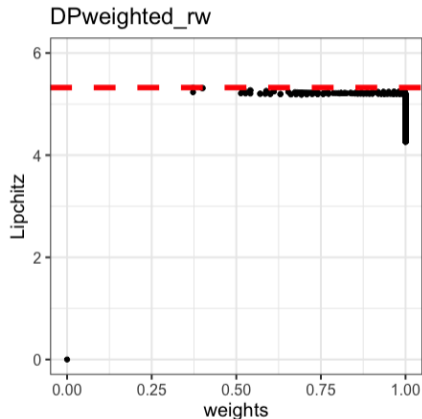
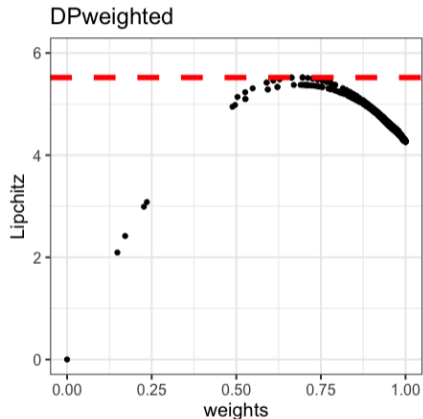
$$\alpha_i^w = k \times \alpha_i \times \frac{\Delta_{\alpha, x}}{\Delta_{\alpha, x_i}} \quad (2)$$

- ▶ Compresses distribution of Δ_{α, x_i}
- ▶ Makes weighting scheme more **efficient**
- ▶ $k \leq 1$ ensures main $\Delta_{\alpha, x}$
- ▶ In the limit the weights are independent Δ_{α, x_i}

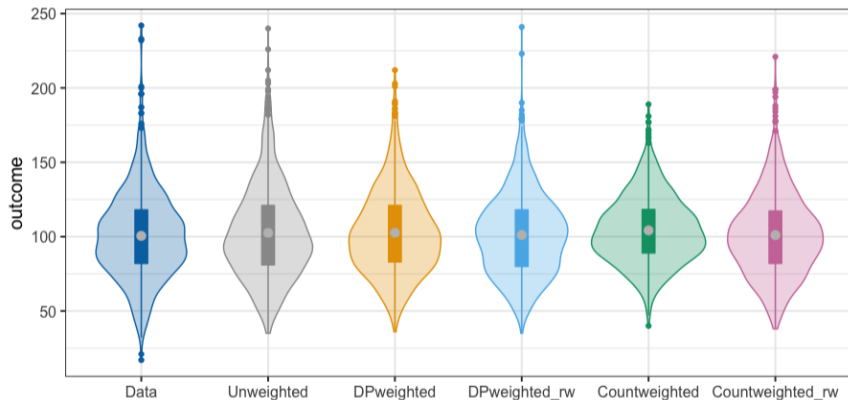
Compression of Δ_{α, x_i} from Reweighting



DP Change in Weighting Efficiency



Improved Estimation of Data Distribution



References I

- Dimitrakakis, C., Nelson, B., Zhang, Z., Mitrokotsa, A. and Rubinstein, B. I. P. (2017), 'Differential privacy for bayesian inference through posterior sampling', *J. Mach. Learn. Res.* **18**(1), 343–381.
URL: <http://dl.acm.org/citation.cfm?id=3122009.3122020>
- Hu, J. and Savitsky, T. D. (2019), 'Bayesian Pseudo Posterior Synthesis for Data Privacy Protection', *arXiv e-prints* p. arXiv:1901.06462.
- Savitsky, T. D., Williams, M. R. and Hu, J. (2019), 'Bayesian Pseudo Posterior Mechanism under Differential Privacy', *arXiv e-prints* p. arXiv:1909.11796.