# General-purpose Statistical Inference with Differential Privacy Guarantees

Zhanyu Wang

Committee: Dr. Jordan Awan (Co-Chair),
Dr. Guang Cheng (Co-Chair),
Dr. Vinayak Rao,
Dr. Christopher W. Clifton.

Department of Statistics, Purdue University

November 17, 2023

**PURDUE**
U N I V E R S I T Y

# Acknowledgement

This thesis is supported in part by

# Table of Contents

# Does Anonymization Guarantee Privacy?

| Data Considered for Sharing | | | | | Voter Registration Records (Identified Resource) | | | |
|---|---|---|---|---|---|---|---|---|
| Age | Zip Code | Gender | Diagnosis | | Birthdate | Zip Code | Gender | Name |
| 15 | 00000 | Male | Diabetes | | 2/2/1989 | 00001 | Female | Alice Smith |
| 21 | 00001 | Female | Influenza | | 3/3/1974 | 10000 | Male | Bob Jones |
| 36 | 10000 | Male | Broken Arm | | 4/4/1919 | 10001 | Female | Charlie Doe |
| 91 | 10001 | Female | Acid Reflux | | | | | |

**Linking two data sources to identity diagnoses.**

Figure: (Department of Health & Human Services) De-identification of sensitive information[1]. Dataset on the left is released without Name. Using another public dataset on the right, we can recover the names in the anonymized dataset.

[1] https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html

Table: Percentage of reconstructed records that exactly agree with the original Census Edited File on location, sex, age, race, and ethnicity[2].

|                                                      | Agreement Rates |
| ---------------------------------------------------- | --------------- |
| Published 2010 Census Tables (swapping)              | 46.5%           |
| Disclosure Avoidance System (differential privacy)   | 15.7%           |

---

[2]https://www2.census.gov/about/partners/cac/sac/meetings/2022-03/
presentation-reconstruction-and-reidentification-of-the-dhc.pdf
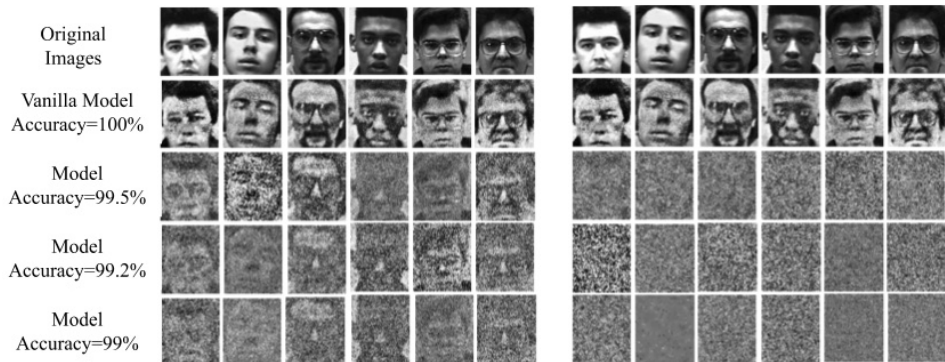
# How to Release a Model Safely?



Figure: Model inversion attack (MIA) results on non-private model trained on Faces94 dataset and differentially privately (DP) trained models (left is record-DP, and right is class-DP)[3].

---

[3]Zhang, Qiuchen, et al. "Broadening differential privacy for deep learning against model inversion attacks." 2020 IEEE International Conference on Big Data (Big Data). IEEE, 2020.

**Statisticians!**

Awesome-Differential-Privacy-for-Statisticians
(my GitHub repo collecting papers on DP+STAT.)

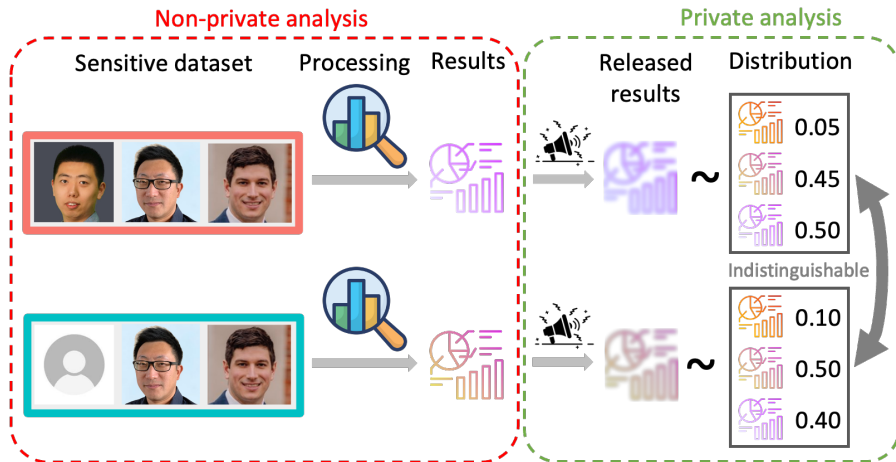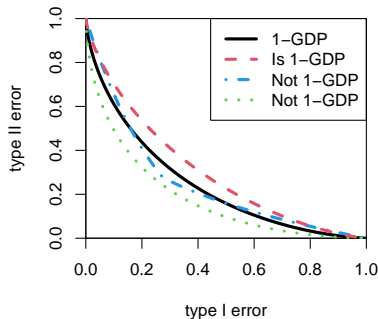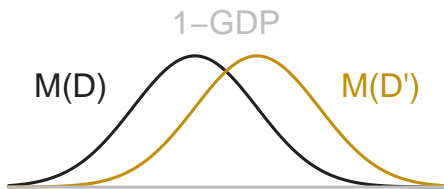# DP: a Probabilistic Measure for Privacy Protection



Figure: The output of the mechanism is roughly the same (approximately indistinguishable) when the input data is slightly changed. This is required for all datasets as input.

▶ A mechanism $M$ is $\mu$-Gaussian DP (Dong, Roth, and Su, 2022, $\mu$-GDP) if for any two datasets $D$, $D'$ differing in one entry, the hypothesis test, using output of $M$, $H_0 : Z \sim M(D), H_1 : Z \sim M(D')$ is never easier than $H_0 : Z \sim N(0,1), H_1 : Z \sim N(\mu, 1)$. (Given type I error, type II is lower bounded.)
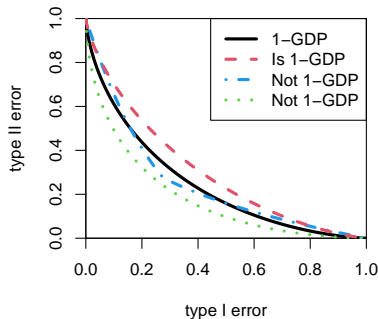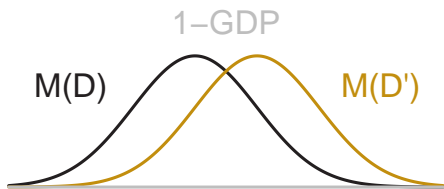
# DP: Formal Definition

▶ A mechanism $M$ is $\mu$-Gaussian DP (Dong, Roth, and Su, 2022, $\mu$-GDP) if for any two datasets $D, D'$ differing in one entry, the hypothesis test, using output of $M$, $H_0 : Z \sim M(D), H_1 : Z \sim M(D')$ is never easier than
$H_0 : Z \sim N(0,1), H_1 : Z \sim N(\mu,1)$. (Given type I error, type II is lower bounded.)



▶ (Our methods also apply to other DP notions like $(\varepsilon, \delta)$-DP or Rényi DP, etc.)

▶ Sensitivity: (largest impact from one individual) The sensitivity of $g(\cdot)$ is

$$\Delta(g) \geq \sup \|g(D) - g(D')\|_2 \, ,$$

where the supremum is over $D$, $D'$ differing in one entry.

▶ **Sensitivity:** (largest impact from one individual) The sensitivity of $g(\cdot)$ is

$$\Delta(g) \geq \sup \|g(D) - g(D')\|_2,$$

where the supremum is over $D$, $D'$ differing in one entry.

▶ **Gaussian Mechanism:** (add noise to protect privacy) If $g$ has sensitivity $\Delta(g)$, then

$$M(D) = g(D) + \xi_{\mathrm{DP}}, \ \xi_{\mathrm{DP}} \sim N\left(0, \left(\frac{\Delta(g)}{\mu}\right)^2\right)$$

satisfies $\mu$-GDP. (Transparency) DP mechanisms are also released for validation.

▶ Sensitivity: (largest impact from one individual) The sensitivity of $g(\cdot)$ is

$$\Delta(g) \geq \sup \|g(D) - g(D')\|_2\,,$$

where the supremum is over $D$, $D'$ differing in one entry.

▶ Gaussian Mechanism: (add noise to protect privacy) If $g$ has sensitivity $\Delta(g)$, then

$$M(D) = g(D) + \xi_{\mathrm{DP}},\ \xi_{\mathrm{DP}} \sim N\left(0, \left(\frac{\Delta(g)}{\mu}\right)^2\right)$$

satisfies $\mu$-GDP. (Transparency) DP mechanisms are also released for validation.

▶ Composition: (more release → less privacy) If $M_1$ and $M_2$ are $\mu_1$-GDP and $\mu_2$-GDP, respectively, then the joint release $(M_1, M_2)$ is $\sqrt{\mu_1^2 + \mu_2^2}$-GDP.

# GDP: Mechanism, Composition, and Post-processing

▶ Sensitivity: (largest impact from one individual) The sensitivity of $g(\cdot)$ is

$$\Delta(g) \geq \sup \|g(D) - g(D')\|_2 \,,$$

where the supremum is over $D$, $D'$ differing in one entry.

▶ Gaussian Mechanism: (add noise to protect privacy) If $g$ has sensitivity $\Delta(g)$, then

$$M(D) = g(D) + \xi_{\mathrm{DP}}, \ \xi_{\mathrm{DP}} \sim N\left(0, \left(\frac{\Delta(g)}{\mu}\right)^2\right)$$

satisfies $\mu$-GDP. (Transparency) DP mechanisms are also released for validation.

▶ Composition: (more release → less privacy) If $M_1$ and $M_2$ are $\mu_1$-GDP and $\mu_2$-GDP, respectively, then the joint release $(M_1, M_2)$ is $\sqrt{\mu_1^2 + \mu_2^2}$-GDP.

▶ Post-processing: (forestall all attackers) If $M(\cdot)$ is $\mu$-GDP, then $\psi(M(\cdot))$ is $\mu$-GDP.

**SOCIAL SCIENCES**

## The use of differential privacy for census data and its impact on redistricting: The case of the 2020 U.S. Census

Christopher T. Kenny[1], Shiro Kuriwaki[2], Cory McCartan[3], Evan T. R. Rosenman[4], Tyler Simko[1], Kosuke Imai[1,3]*



"We find that the [Disclosure Avoidance System] DAS systematically undercounts the population in mixed-race and mixed-partisan precincts, yielding unpredictable racial and partisan biases."

Figure: Mortality risk (relative to current clinical practice) and VKORC1 genotype disclosure risk of DP linear regression used for Warfarin dosing[4].

---

[4]Fredrikson, Matthew, et al. "Privacy in pharmacogenetics: An End-to-End case study of personalized warfarin dosing." 23rd USENIX security symposium. 2014.

**Quantify the uncertainty of the DP output by its sampling distribution.**

Estimate the sampling distribution under DP $\rightarrow$ DP statistical inference.

Focus on frequentist approaches.

**Quantify the uncertainty of the DP output by its sampling distribution.**

Estimate the sampling distribution under DP $\rightarrow$ DP statistical inference.

Focus on frequentist approaches.

- ▶ Model-free. (Part I)
- ▶ Finite-sample valid. (Part II)
- ▶ Optimal. (Part III)
- ▶ Usable when we cannot choose the DP mechanism. (Part II & III)
  - ▶ E.g., post-process the release census data.

**Motivations:**

- ▶ Develop a DP mechanism for <span style="color:red">non-parametric</span> inference.
    - ▶ Build a DP mechanism to enable bootstrap.
- ▶ Perform private inference for <span style="color:red">quantile regression</span>.

---

[5]Wang, Zhanyu, Guang Cheng, and Jordan Awan. "Differentially Private Bootstrap: New Privacy Analysis and Inference Strategies." arXiv:2210.06140 (2022). This work is under review by JMLR.

> Existing privacy guarantees for DP Bootstrap are incorrect, and their confidence intervals have under-coverage.

## DP Bootstrap

▶ Brawner and Honaker (2018); Koskela et al. (2020)
▶ Balle et al. (2018)

---

## DP Parametric Bootstrap
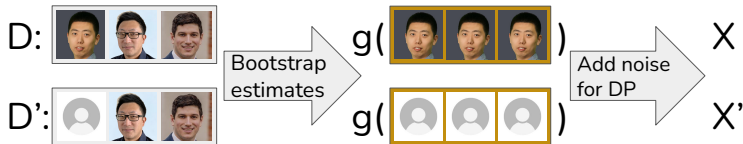
▶ Du et al. (2020); Ferrando et al. (2022); Alabi and Vadhan (2022)

## Bag-of-little bootstrap

▶ Evans et al. (2023); Covington et al. (2021)

▶ The naïve sensitivity is very large $\Rightarrow$ Need to add very large noise for DP.

▶ The naïve sensitivity is very large $\Rightarrow$ Need to add very large noise for DP.



▶ On average, the sensitivity is about the same as without Bootstrap.

# DP Bootstrap Privacy Analysis

▶ The naïve sensitivity is very large $\Rightarrow$ Need to add very large noise for DP.



▶ On average, the sensitivity is about the same as without Bootstrap.

**Theorem: DP Bootstrap Privacy Analysis**

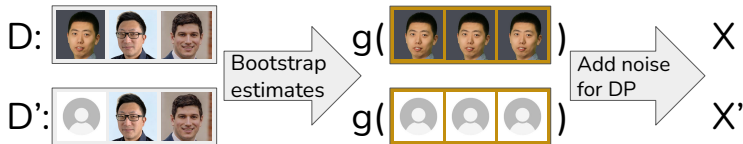▶ If $M$ is $f$-DP, $M \circ \texttt{bootstrap}$ is $f_{\text{boot}}$-DP: $f_{\text{boot}}$ is a tight exact lower bound.

# DP Bootstrap Privacy Analysis

▶ The naïve sensitivity is very large $\Rightarrow$ Need to add very large noise for DP.



▶ On average, the sensitivity is about the same as without Bootstrap.

**Theorem: DP Bootstrap Privacy Analysis**

▶ If $M$ is $f$-DP, $M \circ \texttt{bootstrap}$ is $f_{\text{boot}}$-DP: $f_{\text{boot}}$ is a tight exact lower bound.

▶ If $M$ is $\mu$-GDP, $M \circ \texttt{bootstrap}$ is approximately $\left(\sqrt{2 - 2/e}\right)\mu$-GDP from the above $f_{\text{boot}}$ result when #bootstrap estimates $\to \infty$. ($\sqrt{2 - 2/e} \approx 1.125$)
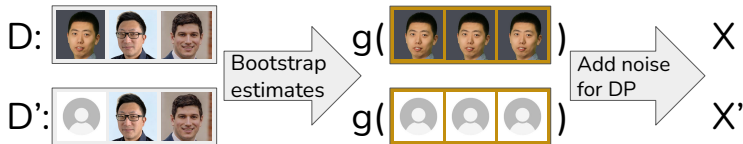
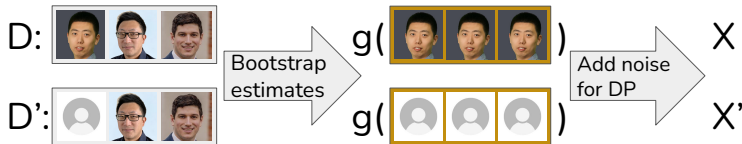▶ The naïve sensitivity is very large $\Rightarrow$ Need to add very large noise for DP.



▶ On average, the sensitivity is about the same as without Bootstrap.

**Theorem: DP Bootstrap Privacy Analysis**

▶ If $M$ is $f$-DP, $M \circ \texttt{bootstrap}$ is $f_{\text{boot}}$-DP: $f_{\text{boot}}$ is a tight exact lower bound.

▶ If $M$ is $\mu$-GDP, $M \circ \texttt{bootstrap}$ is approximately $\left(\sqrt{2 - 2/e}\right)\mu$-GDP from the above $f_{\text{boot}}$ result when #bootstrap estimates $\to \infty$. ($\sqrt{2 - 2/e} \approx 1.125$)

▶ By composition, running $B$ times for $\underline{B \text{ estimates}}$ is $\left(\sqrt{(2 - 2/e)B}\right)\mu$-GDP.

- Sampling distribution is affected by the added noises for DP.

$$M \circ \mathtt{boot}(D) = g(\mathtt{boot}(D)) + \xi_{\mathrm{DP}}$$



- Use deconvolution to recover the distribution of bootstrap estimates from $B$ DP bootstrap estimates and the distribution of added noises.

# Deconvolution for Estimating Sampling Distribution

▶ Sampling distribution is affected by the added noises for DP.

$$M \circ \texttt{boot}(D) = g(\texttt{boot}(D)) + \xi_{\mathrm{DP}}$$



▶ Use deconvolution to recover the distribution of bootstrap estimates from $B$ DP bootstrap estimates and the distribution of added noises.



Sampling distribution estimates

- - - deconvolved private bootstrap
──── non–private bootstrap
······ private bootstrap

$g(D) = \frac{1}{n}\sum_{i=1}^{n} x_i,\ x_i \sim \mathrm{Unif}(0,1),\ n = 10000,\ B = 1000,\ (\sqrt{2 - 2/e})\text{-GDP}.$

# Deconvolution for Estimating Sampling Distribution

▶ Sampling distribution is affected by the added noises for DP.

$$M \circ \texttt{boot}(D) = g(\texttt{boot}(D)) + \xi_{\mathrm{DP}}$$
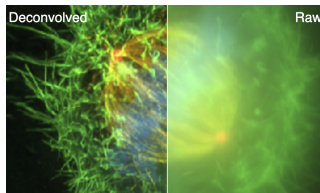


▶ Use deconvolution to recover the distribution of bootstrap estimates from $B$ DP bootstrap estimates and the distribution of added noises.



Sampling distribution estimates
- – – deconvolved private bootstrap
- —— non–private bootstrap
- ⋯⋯ private bootstrap

$g(D) = \frac{1}{n}\sum_{i=1}^{n} x_i,\ x_i \sim \mathrm{Unif}(0,1),\ n = 10000,\ B = 1000,\ (\sqrt{2 - 2/\mathrm{e}})\text{-GDP}.$

▶ Construct DP confidence intervals using quantiles of deconvolved distribution.

# Private Confidence Intervals (CIs) for Quantile Regression

▶ Using the 2016 Canada Census Public Use Microdata, we build 90% CIs for the slope in the <span style="color:red">quantile regression</span> between <span style="color:blue">market income</span> 💰 & <span style="color:blue">shelter cost</span> 🏠.

▶ Using the 2016 Canada Census Public Use Microdata, we build 90% CIs for the slope in the quantile regression between market income 💰 & shelter cost 🏠.

▶ The first DP inference method for quantile regression.

# Private Confidence Intervals (CIs) for Quantile Regression

▶ Using the 2016 Canada Census Public Use Microdata, we build 90% CIs for the slope in the quantile regression between market income 💰 & shelter cost 🏠.

▶ The first DP inference method for quantile regression.

  ▶ For small sample size, DP CIs are a bit wider and more conservative than non-DP;



$\tilde{\theta} = \hat{\theta} + \xi, \ \hat{\theta} = \arg\min_\theta \left( R(\theta) + c\|\theta\|_2^2 \right), \ R(\theta) = \frac{1}{n}\sum_{i=1}^{n}(0.5 - \mathbb{1}(z_i \leq 0))z_i$ where $z_i = y_i - x_i\theta$, $c = 1$, 1-GDP, $\xi \sim N(0, 1/(2n^2))$, $B = 100$.

# Private Confidence Intervals (CIs) for Quantile Regression

▶ Using the 2016 Canada Census Public Use Microdata, we build 90% CIs for the slope in the quantile regression between market income 💰 & shelter cost 🏠.

▶ The first DP inference method for quantile regression.

  ▶ For small sample size, DP CIs are a bit wider and more conservative than non-DP;

  ▶ CIs never contain 0 → <u>significant dependence</u> between 💰 & 🏠.



$\bar{\theta} = \hat{\theta} + \xi, \ \hat{\theta} = \arg\min_\theta \left( R(\theta) + c\|\theta\|_2^2 \right), \ R(\theta) = \frac{1}{n}\sum_{i=1}^n (0.5 - \mathbb{1}(z_i \le 0))z_i$ where $z_i = y_i - x_i\theta$, $c = 1$, 1-GDP, $\xi \sim N(0, 1/(2n^2))$, $B = 100$.

**Contributions:**

1. Propose and analyze a <span style="color:red">non-parametric</span> DP bootstrap framework.
2. The <span style="color:red">first</span> to perform private inference in <span style="color:red">quantile regression</span>.

---

**Limitations:**

1. Bootstrap is an asymptotic method.

**Contributions:**

1. Propose and analyze a non-parametric DP bootstrap framework.
2. The first to perform private inference in quantile regression.

---

**Limitations:**

1. Bootstrap is an asymptotic method.
2. Determining the optimal choice of $B$ is difficult (larger $B \Rightarrow$ more noise & signal).

**Contributions:**

1. Propose and analyze a non-parametric DP bootstrap framework.
2. The first to perform private inference in quantile regression.

---

**Limitations:**

1. Bootstrap is an asymptotic method.
2. Determining the optimal choice of $B$ is difficult (larger $B \Rightarrow$ more noise & signal).
3. Deconvolution is limited to additive noise mechanism (e.g., Gaussian Mechanism).

**Motivations:**

▶ Finite-sample valid coverage/type I errors.

---

[6]Awan, Jordan, and Zhanyu Wang. "Simulation-based, Finite-sample Inference for Privatized Data." arXiv:2303.05328 (2023). This work is under major revision by JASA.

**Motivations:**

► Finite-sample valid coverage/type I errors.

► A general framework that can be used without altering DP mechanisms.

---

[6]Awan, Jordan, and Zhanyu Wang. "Simulation-based, Finite-sample Inference for Privatized Data." arXiv:2303.05328 (2023). This work is under major revision by JASA.

**Motivations:**

▶ Finite-sample valid coverage/type I errors.

▶ A general framework that can be used without altering DP mechanisms.

▶ The privacy mechanism and data generating model are often easy to sample from, enabling simulation-based inference. Our method is inspired by Xie and Wang (2022).

---

[6]Awan, Jordan, and Zhanyu Wang. "Simulation-based, Finite-sample Inference for Privatized Data." arXiv:2305.05328 (2023). This work is under major revision by JASA.

**Example: location-scale normal.** Observe $D := (x_1, \ldots, x_n) \overset{\text{iid}}{\sim} N(\mu^*, \sigma^{*2})$.

▶ Non-private statistic: $\left( m(D) := \frac{1}{n} \sum_{i=1}^{n} x_i, \ \eta^2(D) := \frac{1}{n-1} \sum_{i=1}^{n} \left( x_i - \bar{X} \right)^2 \right)$.

**Example: location-scale normal.** Observe $D := (x_1, \ldots, x_n) \overset{\text{iid}}{\sim} N(\mu^*, \sigma^{*2})$.

▶ Non-private statistic: $\left( m(D) := \frac{1}{n} \sum_{i=1}^{n} x_i, \ \eta^2(D) := \frac{1}{n-1} \sum_{i=1}^{n} \left( x_i - \bar{X} \right)^2 \right)$.

▶ Private statistic, $(\sqrt{2}\varepsilon)$-GDP. $N_1, N_2 \overset{\text{iid}}{\sim} N(0, 1)$, $\mathrm{clamp}_L^U(x) := \max(\min(x, U), L)$.

$$\left( \widetilde{m}(D) := m\left( \mathrm{clamp}_L^U(D) \right) + \frac{U - L}{n\varepsilon} N_1, \ \widetilde{\eta^2}(D) := \eta^2\left( \mathrm{clamp}_L^U(D) \right) + \frac{(U - L)^2}{n\varepsilon} N_2 \right).$$

**Example: location-scale normal.** Observe $D := (x_1, \ldots, x_n) \overset{\text{iid}}{\sim} N(\mu^*, \sigma^{*2})$.

▶ Non-private statistic: $\left( m(D) := \frac{1}{n} \sum_{i=1}^{n} x_i, \ \eta^2(D) := \frac{1}{n-1} \sum_{i=1}^{n} \left( x_i - \bar{X} \right)^2 \right)$.

▶ Private statistic, $(\sqrt{2}\varepsilon)$-GDP. $N_1, N_2 \overset{\text{iid}}{\sim} N(0,1)$, $\text{clamp}_L^U(x) := \max(\min(x, U), L)$.

$$\left( \widetilde{m}(D) := m\left( \text{clamp}_L^U(D) \right) + \frac{U-L}{n\varepsilon} N_1, \ \widetilde{\eta^2}(D) := \eta^2\left( \text{clamp}_L^U(D) \right) + \frac{(U-L)^2}{n\varepsilon} N_2 \right).$$



$\mu^* = 1$, $\sigma^* = 1$, $L = 0$, $U = 3$, $(\sqrt{2})$-GDP.

> Standard techniques are inapplicable or give poor results.

- Likelihood-based inference (Williams and McSherry, 2010)
- Asymptotics (Wang et al., 2018)

**Promising directions:**

- **Parametric bootstrap**
  - Du et al. (2020); Ferrando et al. (2022); Alabi and Vadhan (2022)
- **New asymptotics**
  - Wang et al. (2018, 2019)
- **Bayesian inference via data augmentation MCMC**
  - Ju et al. (2022)

Repro sample (Xie and Wang, 2022) is for likelihood-free simulation-based inference.

# Failure of Naïve Usage of PB under DP

| Privacy guarantee | 1-GDP | 0.5-GDP | 0.3-GDP | 0.1-GDP |
|---|---|---|---|---|
| **Coverage** | 0.803 | 0.806 | 0.804 | 0.819 |

Table: Private 90% confidence intervals by `NOISYVAR+SIM` (Du, Foot, Moniot, Bray, and Groce, 2020) for the population mean of $N(0.5, 1)$. The sample size is 10000.

# Failure of Naïve Usage of PB under DP

| Privacy guarantee | 1-GDP | 0.5-GDP | 0.3-GDP | 0.1-GDP |
|---|---|---|---|---|
| Coverage | 0.803 | 0.806 | 0.804 | 0.819 |

Table: Private 90% confidence intervals by `NOISYVAR+SIM` (Du, Foot, Moniot, Bray, and Groce, 2020) for the population mean of $N(0.5, 1)$. The sample size is 10000.

| Sample size | 100 | 200 | 500 | 1000 | 2000 | 5000 |
|---|---|---|---|---|---|---|
| Type I error | 0.017 | 0.045 | 0.118 | 0.186 | 0.361 | 0.674 |

Table: Private hypothesis testings (level 0.05) using DP Monte Carlo tests (Alabi and Vadhan, 2022) on $H_0 : \beta_1^* = 0$ and $H_1 : \beta_1^* \neq 0$ with a regression model $Y = \beta_0^* + X\beta_1^* + \epsilon$ under 1-GDP.

Example (Nonprivate Location Normal)

## Example (Nonprivate Location Normal)

▶ $x_1, \ldots, x_n \overset{\text{iid}}{\sim} N(\theta^*, 1)$. Observe the statistic $s^* := \overline{x}$.

## Example (Nonprivate Location Normal)

- $x_1, \ldots, x_n \overset{\text{iid}}{\sim} N(\theta^*, 1)$. Observe the statistic $s^* := \overline{x}$.
- A $(1 - \alpha)$-prediction interval for $s^*$ is $B_\alpha(\theta^*) = \left[\theta^* - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \theta^* + \frac{z_{1-\alpha/2}}{\sqrt{n}}\right]$.

## Example (Nonprivate Location Normal)

- $x_1, \ldots, x_n \overset{\text{iid}}{\sim} N(\theta^*, 1)$. Observe the statistic $s^* := \overline{x}$.

- A $(1 - \alpha)$-prediction interval for $s^*$ is $B_\alpha(\theta^*) = \left[ \theta^* - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \theta^* + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right]$.

- The confidence interval is

$$\Gamma_\alpha(s^*) = \{\theta \mid s^* \in B_\alpha(\theta)\}$$
$$= \left[ \overline{x} - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \overline{x} + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right].$$

## Example (Nonprivate Location Normal)

▶ $x_1, \ldots, x_n \overset{\text{iid}}{\sim} N(\theta^*, 1)$. Observe the statistic $s^* := \overline{x}$.

▶ A $(1 - \alpha)$-prediction interval for $s^*$ is $B_\alpha(\theta^*) = \left[ \theta^* - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \theta^* + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right]$.

▶ The confidence interval is

$$\Gamma_\alpha(s^*) = \{\theta \mid s^* \in B_\alpha(\theta)\}$$
$$= \left[ \overline{x} - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \overline{x} + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right].$$

## Example (Nonprivate Location Normal)

- $x_1, \ldots, x_n \overset{\text{iid}}{\sim} N(\theta^*, 1)$. Observe the statistic $s^* := \overline{x}$.
- A $(1-\alpha)$-prediction interval for $s^*$ is $B_\alpha(\theta^*) = \left[ \theta^* - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \theta^* + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right]$.

- The confidence interval is

$$\Gamma_\alpha(s^*) = \{\theta \mid s^* \in B_\alpha(\theta)\}$$
$$= \left[ \overline{x} - \frac{z_{1-\alpha/2}}{\sqrt{n}}, \ \overline{x} + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right].$$



A confidence set $\Gamma_\alpha(s^*)$ can be constructed by inverting a prediction set $B_\alpha(\theta)$.

**Data-generating model:**

$$D^* := G_{\mathrm{data}}(\theta^*, u_{\mathrm{data}}).$$

▶ $\theta^*$ is unknown, $u_{\mathrm{data}} \sim F_{\mathrm{data}}$ is a random seed, $G_{\mathrm{data}}$ is a deterministic function.

**Data-generating model:**

$$D^* := G_{\mathrm{data}}(\theta^*, u_{\mathrm{data}}).$$

- $\theta^*$ is unknown, $u_{\mathrm{data}} \sim F_{\mathrm{data}}$ is a random seed, $G_{\mathrm{data}}$ is a deterministic function.
- $D^* := (x_1, \ldots, x_n) \overset{\mathrm{iid}}{\sim} N(\mu^*, \sigma^{*2}) \Leftrightarrow D^* := \mu^* + \sigma^* u, \ u \sim N(0, I_{n \times n})$.

**Data-generating model:**
$$D^* := G_{\mathrm{data}}(\theta^*, u_{\mathrm{data}}).$$

▶ $\theta^*$ is unknown, $u_{\mathrm{data}} \sim F_{\mathrm{data}}$ is a random seed, $G_{\mathrm{data}}$ is a deterministic function.

▶ $D^* := (x_1, \ldots, x_n) \overset{\mathrm{iid}}{\sim} N(\mu^*, \sigma^{*2}) \Leftrightarrow D^* := \mu^* + \sigma^* u, \ u \sim N(0, I_{n \times n})$.

**Private statistics:**
$$s^* := G_{\mathrm{privacy}}(D^*, u_{\mathrm{privacy}}).$$

▶ DP mechanism $G_{\mathrm{privacy}}$ contains extra uncertainty $u_{\mathrm{privacy}} \sim F_{\mathrm{privacy}}$,

**Data-generating model:**

$$D^* := G_{\mathrm{data}}(\theta^*, u_{\mathrm{data}}).$$

▶ $\theta^*$ is unknown, $u_{\mathrm{data}} \sim F_{\mathrm{data}}$ is a random seed, $G_{\mathrm{data}}$ is a deterministic function.

▶ $D^* := (x_1, \ldots, x_n) \overset{\mathrm{iid}}{\sim} N(\mu^*, \sigma^{*2}) \Leftrightarrow D^* := \mu^* + \sigma^* u, \ u \sim N(0, I_{n \times n})$.

**Private statistics:**

$$s^* := G_{\mathrm{privacy}}(D^*, u_{\mathrm{privacy}}).$$

▶ DP mechanism $G_{\mathrm{privacy}}$ contains extra uncertainty $u_{\mathrm{privacy}} \sim F_{\mathrm{privacy}}$,

▶ e.g., Gaussian Mechanism: $s^* := g(D^*) + u_{\mathrm{privacy}}$.

**Data-generating model:**

$$D^* := G_{\text{data}}(\theta^*, u_{\text{data}}).$$

- ▶ $\theta^*$ is unknown, $u_{\text{data}} \sim F_{\text{data}}$ is a random seed, $G_{\text{data}}$ is a deterministic function.
- ▶ $D^* := (x_1, \ldots, x_n) \overset{\text{iid}}{\sim} N(\mu^*, \sigma^{*2}) \Leftrightarrow D^* := \mu^* + \sigma^* u, \ u \sim N(0, I_{n \times n})$.

**Private statistics:**

$$s^* := G_{\text{privacy}}(D^*, u_{\text{privacy}}).$$

- ▶ DP mechanism $G_{\text{privacy}}$ contains extra uncertainty $u_{\text{privacy}} \sim F_{\text{privacy}}$,
- ▶ e.g., Gaussian Mechanism: $s^* := g(D^*) + u_{\text{privacy}}$.

Combine them and write the **generating equation** as $s^* \overset{d}{=} G(\theta^*, u)$.

**Data-generating model:**
$$D^* := G_{\mathrm{data}}(\theta^*, u_{\mathrm{data}}).$$

▶ $\theta^*$ is unknown, $u_{\mathrm{data}} \sim F_{\mathrm{data}}$ is a random seed, $G_{\mathrm{data}}$ is a deterministic function.

▶ $D^* := (x_1, \ldots, x_n) \overset{\mathrm{iid}}{\sim} N(\mu^*, \sigma^{*2}) \Leftrightarrow D^* := \mu^* + \sigma^* u, \ u \sim N(0, I_{n \times n})$.

**Private statistics:**
$$s^* := G_{\mathrm{privacy}}(D^*, u_{\mathrm{privacy}}).$$

▶ DP mechanism $G_{\mathrm{privacy}}$ contains extra uncertainty $u_{\mathrm{privacy}} \sim F_{\mathrm{privacy}}$,

▶ e.g., Gaussian Mechanism: $s^* := g(D^*) + u_{\mathrm{privacy}}$.

Combine them and write the **generating equation** as $s^* \overset{d}{=} G(\theta^*, u)$.

This setup (and our method) applies to all settings with low-d summary statistics $s^*$.

▶ Simulate new random seeds $u_i \stackrel{\text{iid}}{\sim} P$ and fix them. Let $s_i(\theta) = G(\theta, u_i)$.

▶ Simulate new random seeds $u_i \overset{\text{iid}}{\sim} P$ and fix them. Let $s_i(\theta) = G(\theta, u_i)$.

▶ Let $\mathbf{S}(\theta) = \{s^*, s_1(\theta), \ldots, s_R(\theta)\}$. Then, $\mathbf{S}(\theta^*)$ is a set of exchangeable r.v.s.

▶ Simulate new random seeds $u_i \overset{\text{iid}}{\sim} P$ and fix them. Let $s_i(\theta) = G(\theta, u_i)$.

▶ Let $\mathbf{S}(\theta) = \{s^*, s_1(\theta), \ldots, s_R(\theta)\}$. Then, $\mathbf{S}(\theta^*)$ is a set of exchangeable r.v.s.

▶ $s_i(\theta) \approx s^* \Rightarrow \theta \approx \theta^*$. Define $T_{\mathbf{S}}(s)$ as the "closeness" between $s$ and $\mathbf{S}$.

# Simulation-Based Confidence Sets by Repro Samples

▶ Simulate new random seeds $u_i \overset{\text{iid}}{\sim} P$ and fix them. Let $s_i(\theta) = G(\theta, u_i)$.

▶ Let $\mathbf{S}(\theta) = \{s^*, s_1(\theta), \ldots, s_R(\theta)\}$. Then, $\mathbf{S}(\theta^*)$ is a set of exchangeable r.v.s.

▶ $s_i(\theta) \approx s^* \Rightarrow \theta \approx \theta^*$. Define $T_{\mathbf{S}}(s)$ as the "closeness" between $s$ and $\mathbf{S}$.

▶ $T_{\mathbf{S}}(s^*), T_{\mathbf{S}}(s_1(\theta^*)), \ldots, T_{\mathbf{S}}(s_R(\theta^*))$
exchangeable w/ order statistics $T_{(i)}^{\theta^*}$.

$$\mathbb{P}\left( T_{\mathbf{S}(\theta^*)}(s^*) \in \left[ T_{(\alpha(R+1)+1)}^{\theta^*}, T_{(R+1)}^{\theta^*} \right] \right) \geq 1-\alpha.$$



larger $T \Rightarrow s_i(\theta)$ closer to $s^* \Rightarrow$ better $\theta$

# Simulation-Based Confidence Sets by Repro Samples

▶ Simulate new random seeds $u_i \overset{iid}{\sim} P$ and fix them. Let $s_i(\theta) = G(\theta, u_i)$.

▶ Let $\mathbf{S}(\theta) = \{s^*, s_1(\theta), \ldots, s_R(\theta)\}$. Then, $\mathbf{S}(\theta^*)$ is a set of exchangeable r.v.s.
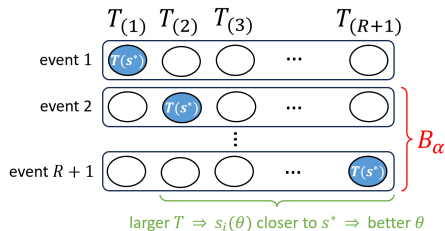
▶ $s_i(\theta) \approx s^* \Rightarrow \theta \approx \theta^*$. Define $T_{\mathbf{S}}(s)$ as the "closeness" between $s$ and $\mathbf{S}$.

▶ $T_{\mathbf{S}}(s^*), T_{\mathbf{S}}(s_1(\theta^*)), \ldots, T_{\mathbf{S}}(s_R(\theta^*))$
  exchangeable w/ order statistics $T_{(i)}^{\theta^*}$.

  $$\mathbb{P}\left( T_{\mathbf{S}(\theta^*)}(s^*) \in \left[ T_{(\alpha(R+1)+1)}^{\theta^*}, T_{(R+1)}^{\theta^*} \right] \right) \geq 1-\alpha.$$

▶ For general $\theta$, define $B_\alpha(\theta)$ as
  $\left\{ T_{\mathbf{S}(\theta)}(s^*) \in \left[ T_{(\alpha(R+1)+1)}^{\theta}, T_{(R+1)}^{\theta} \right] \right\}$.



larger $T \Rightarrow s_i(\theta)$ closer to $s^* \Rightarrow$ better $\theta$

# Simulation-Based Confidence Sets by Repro Samples

▶ Simulate new random seeds $u_i \overset{\text{iid}}{\sim} P$ and fix them. Let $s_i(\theta) = G(\theta, u_i)$.

▶ Let $\mathbf{S}(\theta) = \{s^*, s_1(\theta), \ldots, s_R(\theta)\}$. Then, $\mathbf{S}(\theta^*)$ is a set of exchangeable r.v.s.

▶ $s_i(\theta) \approx s^* \Rightarrow \theta \approx \theta^*$. Define $T_{\mathbf{S}}(s)$ as the "closeness" between $s$ and $\mathbf{S}$.

▶ $T_{\mathbf{S}}(s^*), T_{\mathbf{S}}(s_1(\theta^*)), \ldots, T_{\mathbf{S}}(s_R(\theta^*))$
exchangeable w/ order statistics $T_{(i)}^{\theta^*}$.

$$\mathbb{P}\left(T_{\mathbf{S}(\theta^*)}(s^*) \in \left[T_{(\alpha(R+1)+1)}^{\theta^*}, T_{(R+1)}^{\theta^*}\right]\right) \geq 1-\alpha.$$

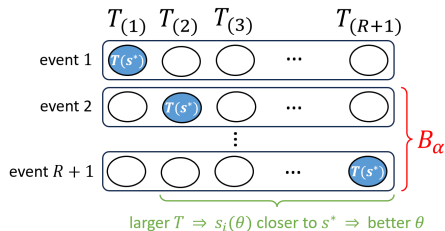▶ For general $\theta$, define $B_\alpha(\theta)$ as
$\left\{T_{\mathbf{S}(\theta)}(s^*) \in \left[T_{(\alpha(R+1)+1)}^{\theta}, T_{(R+1)}^{\theta}\right]\right\}$.

▶ $\Gamma_\alpha := \{\theta \mid \mathbb{1}(B_\alpha(\theta)) = 1\}$ is a
$(1-\alpha)$-confidence set for $\theta^*$.



larger $T \Rightarrow s_i(\theta)$ closer to $s^* \Rightarrow$ better $\theta$

## Theorem: Confidence set from simulated (repro) samples

Set $\mathbf{S} = (s^*, s_1(\theta), \ldots, s_R(\theta))$ and $\left\{ T_{(i)}^{\theta} \right\}_{i=1}^{R+1}$ be order statistics of

$$T(s^*; \mathbf{S}), T(s_1(\theta); \mathbf{S}), \ldots, T(s_R(\theta); \mathbf{S}),$$

where $T$ is permutation-invariant in $\mathbf{S}$. Then $\left\{ T_{(i)}^{\theta^*} \right\}_{i=1}^{R+1}$ are exchangeable. If lower values of $T$ indicate unusual data points, then, a $(1 - \alpha)$-confidence set is

$$\Gamma_\alpha(s^*, u) := \left\{ \theta \mid T(s^*; \mathbf{S}) \in \left[ T_{(\lfloor \alpha(R+1) \rfloor + 1)}^{\theta}, T_{(R+1)}^{\theta} \right] \right\}.$$

**Key insights:**

1. Include $s^*$ in $\mathbf{S}$ to ensure exchangeability from permutation-invariance.
2. Prediction set from order statistics, like conformal prediction (Vovk et al., 2005).

▶ Most statistical depths are permutation-invariant, and unusual points have lower depth, e.g., Mahalanobis depth: $T(s; \mathbf{S}) = \left[1 + (s - \mu_{\mathbf{S}})^{\mathsf{T}} \Sigma_{\mathbf{S}}^{-1}(s - \mu_{\mathbf{S}})\right]^{-1}$, where $(\mu_{\mathbf{S}}, \Sigma_{\mathbf{S}})$ is sample (mean, covariance) of $\mathbf{S}$.

▶ Comparing different depths with $s^* := \left(\widetilde{m}(D), \, \widetilde{\eta^2}(D)\right) = (1, 0.75)$.

▶ We can leverage exchangeability to derive $p$-values as well.

**Theorem: Hypothesis testing $p$-value**

If $T$ is a depth function taking value in $(0, 1)$, then

$$p = \frac{1}{R+1} \left\lfloor \sup_{\theta \in \Theta_0} \left[ \# \left\{ i \mid T_{(i)}^{\theta} \leq T(s^*; \mathbf{S}) \right\} + T(s^*; \mathbf{S}) \right] \right\rfloor$$

is a valid $p$-value for $H_0 : \theta^* \in \Theta_0$.

The main competitor for general frequentist inference for privatized data is the parametric bootstrap (PB).[7]

▶ Du et al. (2020); Ferrando et al. (2022); Alabi and Vadhan (2022).



---

[7]Figure credit: Boos, Dennis, and Leonard Stefanski. "Efron's bootstrap." Significance (2010).

The main competitor for general frequentist inference for privatized data is the parametric bootstrap (PB).[7]

▶ Du et al. (2020); Ferrando et al. (2022); Alabi and Vadhan (2022).

PB uses a parametric model for inference, while Repro uses a data-generating equation. However,

▶ With a biased estimator, PB can give inaccurate inferences.



---

[7]Figure credit: Boos, Dennis, and Leonard Stefanski. "Efron's bootstrap." Significance (2010).

The main competitor for general frequentist inference for privatized data is the parametric bootstrap (PB).[7]

▶ Du et al. (2020); Ferrando et al. (2022); Alabi and Vadhan (2022).

PB uses a parametric model for inference, while Repro uses a data-generating equation. However,

▶ With a biased estimator, PB can give inaccurate inferences.

▶ PB lacks finite sample guarantees.



---

[7]Figure credit: Boos, Dennis, and Leonard Stefanski. "Efron's bootstrap." Significance (2010).

# Location-Scale Normal

- Suppose that $D := (x_1, \ldots, x_n)$, $x_i \overset{\text{iid}}{\sim} N(\mu^*, \sigma^*)$. Build CIs for $\mu^*$ and $\sigma^*$.
- True parameter $\theta^* := (\mu^*, \sigma^*)$.
- DP statistic $s := \left( \widetilde{m}, \widetilde{\eta^2} \right)$ satisfies $(\sqrt{2}\varepsilon)$-GDP.

$\widetilde{m}(D) := m \left( \text{clamp}_L^U(D) \right) + \frac{U-L}{n\varepsilon} N_1$, $\widetilde{\eta^2}(D) := \eta^2 \left( \text{clamp}_L^U(D) \right) + \frac{(U-L)^2}{n\varepsilon} N_2$, where $N_1, N_2 \overset{\text{iid}}{\sim} N(0, 1)$, $\text{clamp}_L^U(x) := \max(\min(x, U), L)$.

| Method (95% CI) | Coverage | | Average width | |
|---|---|---|---|---|
| | $\mu^*$ | $\sigma^*$ | $\mu^*$ | $\sigma^*$ |
| Repro Sample | 0.989 (0.003) | 0.998 (0.001) | 0.599 (0.003) | 0.758 (0.005) |
| Parametric Bootstrap | 0.688 (0.015) | 0.003 (0.001) | 0.311 (0.001) | 0.291 (0.001) |

$\mu^* = 1$, $\sigma^* = 1$, $L = 0$, $U = 3$, $R = 200$, $(\sqrt{2})$-GDP.

# Hypothesis Testing for Linear Regression

- This problem setting is used by (Alabi and Vadhan, 2022).
- Test $H_0 : \beta_1^* = 0$ and $H_1 : \beta_1^* \neq 0$ with $Y = \beta_0^* + X\beta_1^* + \epsilon$.
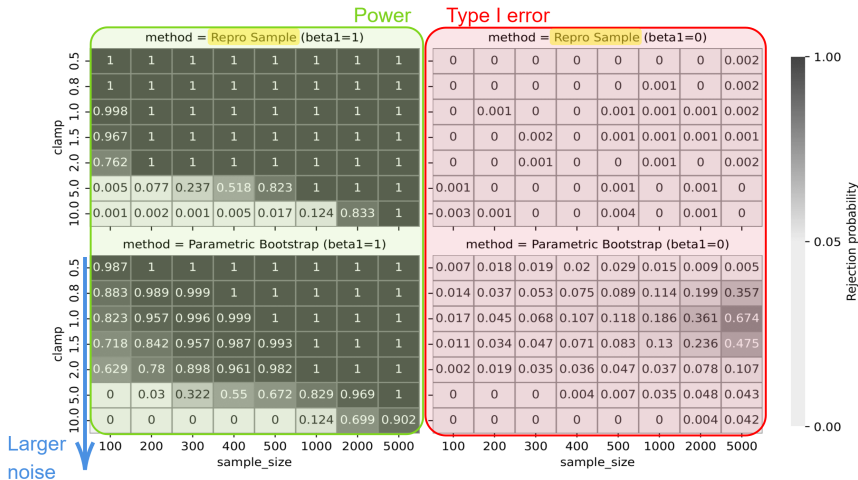
# Hypothesis Testing for Linear Regression

▶ This problem setting is used by (Alabi and Vadhan, 2022).

▶ Test $H_0 : \beta_1^* = 0$ and $H_1 : \beta_1^* \neq 0$ with $Y = \beta_0^* + X\beta_1^* + \epsilon$.

▶ True parameter $\theta^* := (\beta_1^*, \beta_0^*, \mathbb{E}[X], \mathrm{Var}(X), \mathrm{Var}(\varepsilon))$.

▶ DP statistic $s := \left( \tilde{\tilde{x}}, \widetilde{x^2}, \tilde{\tilde{y}}, \widetilde{xy}, \widetilde{y^2} \right)$ satisfies $\mu$-GDP:

Set clamp parameter to be $\Delta$, $[x_i]_{-\Delta}^{\Delta} := \max(\min(x, \Delta), -\Delta)$.

$$\tilde{\tilde{x}} := \frac{1}{n} \sum_{i=1}^{n} [x_i]_{-\Delta}^{\Delta} + \frac{2\Delta}{(\mu/\sqrt{5})n} N_1, \qquad \widetilde{x^2} := \frac{1}{n} \sum_{i=1}^{n} [x_i^2]_0^{\Delta^2} + \frac{\Delta^2}{(\mu/\sqrt{5})n} N_2,$$

$$\tilde{\tilde{y}} := \frac{1}{n} \sum_{i=1}^{n} [y_i]_{-\Delta}^{\Delta} + \frac{2\Delta}{(\mu/\sqrt{5})n} N_3, \qquad \widetilde{xy} := \frac{1}{n} \sum_{i=1}^{n} [x_i y_i]_{-\Delta^2}^{\Delta^2} + \frac{2\Delta^2}{(\mu/\sqrt{5})n} N_4,$$

$$\widetilde{y^2} := \frac{1}{n} \sum_{i=1}^{n} [y_i^2]_0^{\Delta^2} + \frac{\Delta^2}{(\mu/\sqrt{5})n} N_5, \qquad \text{where} \quad N_i \overset{\text{iid}}{\sim} N(0, 1).$$

▶ Compare rejection probabilities (level 0.05) to PB (Alabi and Vadhan, 2022).

**Contributions:** Expand Repro for finite-sample inference for privatized data.

1. We ensure valid coverage/type I errors, even accounting for Monte Carlo errors;
2. Our method is post-processor and can be applied without additional DP budget.
3. We apply it to many private inference problems and compare it to other methods.

---

**Limitations:**

1. Confidence set may be discontinuous.

**Contributions:** Expand Repro for finite-sample inference for privatized data.

1. We ensure valid coverage/type I errors, even accounting for Monte Carlo errors;
2. Our method is post-processor and can be applied without additional DP budget.
3. We apply it to many private inference problems and compare it to other methods.

---

**Limitations:**

1. Confidence set may be discontinuous.
2. Often conservative.

**Contributions:** Expand Repro for finite-sample inference for privatized data.

1. We ensure valid coverage/type I errors, even accounting for Monte Carlo errors;
2. Our method is post-processor and can be applied without additional DP budget.
3. We apply it to many private inference problems and compare it to other methods.

---

**Limitations:**

1. Confidence set may be discontinuous.
2. Often conservative.
   - ▶ Using pivotal summary statistics gives better performance, e.g., in logistic regression.

**Contributions:** Expand Repro for finite-sample inference for privatized data.

1. We ensure valid coverage/type I errors, even accounting for Monte Carlo errors;
2. Our method is post-processor and can be applied without additional DP budget.
3. We apply it to many private inference problems and compare it to other methods.

---

**Limitations:**

1. Confidence set may be discontinuous.
2. Often conservative.
   ▶ Using pivotal summary statistics gives better performance, e.g., in logistic regression.
3. Optimizing over the nuisance parameters is computationally expensive.

**Motivations:**

▶ Repro method is over-conservative and has no optimality guarantee.

▶ Existing parametric bootstrap (PB) gives biased results due to clamping.

▶ Need an estimator for PB: consistent & achieving the optimal asymptotic variance.

---

[8]Wang, Zhanyu, and Jordan Awan. "Debiased Parametric Bootstrap Inference on Privatized Data."
This work is presented in TPDP 2023 and under preparation for journal submission.

# Bias in Naïve DP Estimates due to Clamping

**Example: location-scale normal.** Observe $D := (x_1, \ldots, x_n) \overset{\text{iid}}{\sim} N(\mu^*, \sigma^{*2})$.

▶ Non-private statistic: $\left( m(D) := \frac{1}{n} \sum_{i=1}^{n} x_i, \ \eta^2(D) := \frac{1}{n-1} \sum_{i=1}^{n} \left( x_i - \bar{X} \right)^2 \right)$.

▶ Private statistic, $(\sqrt{2}\varepsilon)$-GDP. $N_1, N_2 \overset{\text{iid}}{\sim} N(0, 1)$, $\text{clamp}_L^U(x) := \max(\min(x, U), L)$.

$$\left( \widetilde{m}(D) := m\left(\text{clamp}_L^U(D)\right) + \frac{U - L}{n\varepsilon} N_1, \ \widetilde{\eta^2}(D) := \eta^2\left(\text{clamp}_L^U(D)\right) + \frac{(U-L)^2}{n\varepsilon} N_2 \right).$$



$\mu^* = 1, \sigma^* = 1, L = 0, U = 3, (\sqrt{2})$-GDP.

- $D \sim F(x|\theta^*)$ with distribution $F$, true parameter $\theta^*$.
  Compute summary statistic $s(D)$.

# Parametric Bootstrap (PB)

▶ $D \sim F(x|\theta^*)$ with distribution $F$,
true parameter $\theta^*$.
Compute summary statistic $s(D)$.

▶ Parameter of interest $\tau^* := \tau(\theta^*)$.

# Parametric Bootstrap (PB)

▶ $D \sim F(x|\theta^*)$ with distribution $F$, true parameter $\theta^*$.
Compute summary statistic $s(D)$.

▶ Parameter of interest $\tau^* := \tau(\theta^*)$.

▶ Estimate $\theta^*$, $\tau^*$ by $\hat{\theta}(s)$, $\hat{\tau}(s)$.

- $D \sim F(x|\theta^*)$ with distribution $F$, true parameter $\theta^*$.
  Compute summary statistic $s(D)$.
- Parameter of interest $\tau^* := \tau(\theta^*)$.
- Estimate $\theta^*$, $\tau^*$ by $\hat{\theta}(s)$, $\hat{\tau}(s)$.
- Let $D_b \sim F(x|\hat{\theta}(s))$, $s_b = s(D_b)$. The PB estimator of $\tau^*$ is $\hat{\tau}(s_b)$.

- $D \sim F(x|\theta^*)$ with distribution $F$, true parameter $\theta^*$.
  Compute summary statistic $s(D)$.
- Parameter of interest $\tau^* := \tau(\theta^*)$.
- Estimate $\theta^*$, $\tau^*$ by $\hat{\theta}(s)$, $\hat{\tau}(s)$.
- Let $D_b \sim F(x|\hat{\theta}(s))$, $s_b = s(D_b)$.
  The PB estimator of $\tau^*$ is $\hat{\tau}(s_b)$.



- Distributions of $\sqrt{n}(\hat{\tau}(D) - \tau^*)$ and $\sqrt{n}(\hat{\tau}(D_b) - \hat{\tau}(D))$ are $H_n(\theta^*)$ and $H_n(\hat{\theta}(D))$.

- $D \sim F(x|\theta^*)$ with distribution $F$, true parameter $\theta^*$.
  Compute summary statistic $s(D)$.

- Parameter of interest $\tau^* := \tau(\theta^*)$.

- Estimate $\theta^*$, $\tau^*$ by $\hat{\theta}(s)$, $\hat{\tau}(s)$.

- Let $D_b \sim F(x|\hat{\theta}(s))$, $s_b = s(D_b)$.
  The PB estimator of $\tau^*$ is $\hat{\tau}(s_b)$.



- Distributions of $\sqrt{n}(\hat{\tau}(D) - \tau^*)$ and $\sqrt{n}(\hat{\tau}(D_b) - \hat{\tau}(D))$ are $H_n(\theta^*)$ and $H_n(\hat{\theta}(D))$.

- PB consistency: $H_n(\hat{\theta}(D)) \xrightarrow{P} H_n(\theta^*)$. $\Rightarrow$ Asymptotically valid CIs & HTs by PB.
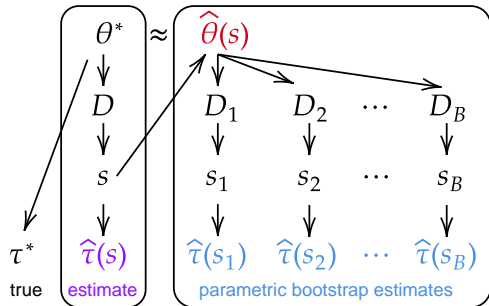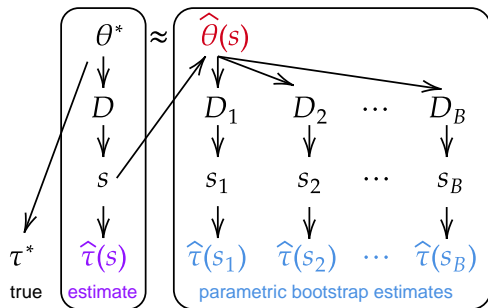
# Parametric Bootstrap (PB)

- $D \sim F(x|\theta^*)$ with distribution $F$, true parameter $\theta^*$.
  Compute summary statistic $s(D)$.
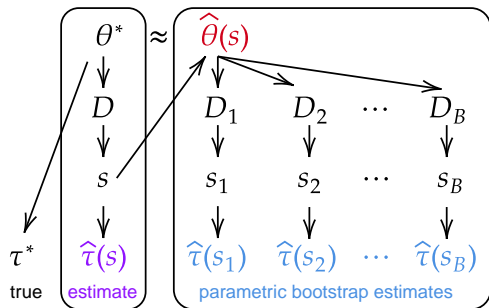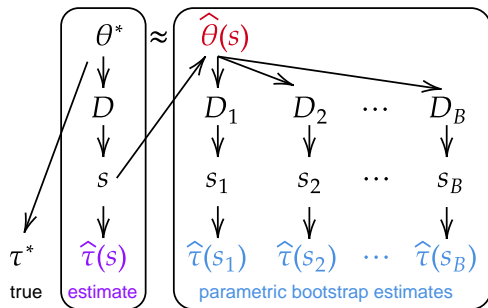- Parameter of interest $\tau^* := \tau(\theta^*)$.
- Estimate $\theta^*$, $\tau^*$ by $\hat{\theta}(s)$, $\hat{\tau}(s)$.
- Let $D_b \sim F(x|\hat{\theta}(s))$, $s_b = s(D_b)$.
  The PB estimator of $\tau^*$ is $\hat{\tau}(s_b)$.



- Distributions of $\sqrt{n}(\hat{\tau}(D) - \tau^*)$ and $\sqrt{n}(\hat{\tau}(D_b) - \hat{\tau}(D))$ are $H_n(\theta^*)$ and $H_n(\hat{\theta}(D))$.
- PB consistency: $H_n(\hat{\theta}(D)) \xrightarrow{P} H_n(\theta^*)$. $\Rightarrow$ Asymptotically valid CIs & HTs by PB.

> When $\hat{\theta}(D)$ in PB is biased, the consistency may not hold.

> Naïve usage of PB under DP gives biased results.

**Parametric bootstrap**

▶ Du et al. (2020); Ferrando et al. (2022); Alabi and Vadhan (2022)

**Bag-of-little bootstrap**

▶ Evans et al. (2023); Covington et al. (2021)

**Indirect inference for bias correction**

▶ Gourieroux et al. (1993); Jiang and Turnbull (2004); Guerrier et al. (2019)

> The privacy mechanism and data generating model are often **easy to sample from**, enabling simulation-based inference.

▶ Write the **generating equation** as $s^* \overset{d}{=} G(\theta^*, u)$.

> The privacy mechanism and data generating model are often **easy to sample** from, enabling simulation-based inference.

▶ Write the **generating equation** as $s^* \stackrel{d}{=} G(\theta^*, u)$.

▶ Fix the randomness $\{u_i\}_{i=1}^{R}$ in generating $s_i(\theta) = G(\theta, u_i)$.

> The privacy mechanism and data generating model are often **easy to sample** from, enabling simulation-based inference.

▶ Write the **generating equation** as $s^* \stackrel{d}{=} G(\theta^*, u)$.

▶ Fix the randomness $\{u_i\}_{i=1}^R$ in generating $s_i(\theta) = G(\theta, u_i)$.

▶ If $s_i(\theta)$ is close to $s^*$, $\theta$ is close to $\theta^*$.

> The privacy mechanism and data generating model are often **easy to sample** from, enabling simulation-based inference.

- Write the **generating equation** as $s^* \overset{d}{=} G(\theta^*, u)$.
- Fix the randomness $\{u_i\}_{i=1}^{R}$ in generating $s_i(\theta) = G(\theta, u_i)$.
- If $s_i(\theta)$ is close to $s^*$, $\theta$ is close to $\theta^*$.
- Find the $\theta$ generating $s_i(\theta)$ closest to $s^*$. Use $\|x\|_{\Omega} := \sqrt{x^{\intercal} \Omega x}$ as a metric.

# Indirect Estimator (Gourieroux, Monfort, and Renault, 1993)

> The privacy mechanism and data generating model are often **easy to sample** from, enabling simulation-based inference.

- Write the **generating equation** as $s^* \overset{d}{=} G(\theta^*, u)$.
- Fix the randomness $\{u_i\}_{i=1}^R$ in generating $s_i(\theta) = G(\theta, u_i)$.
- If $s_i(\theta)$ is close to $s^*$, $\theta$ is close to $\theta^*$.
- Find the $\theta$ generating $s_i(\theta)$ closest to $s^*$. Use $\|x\|_\Omega := \sqrt{x^\mathsf{T} \Omega x}$ as a metric.

Definition (Indirect estimator)

$$\hat{\theta}_{\mathrm{IND}} := \operatorname*{arg\,min}_{\theta \in \Theta} \left\| s^* - \frac{1}{R} \sum_{i=1}^R s_i(\theta) \right\|_\Omega.$$

**Theorem: Indirect Estimator Consistency**

1. (Gourieroux et al., 1993) $\hat{\theta}_{\mathrm{IND}}$ is a consistent estimator of $\theta^*$, and $\sqrt{n}\left(\hat{\theta}_{\mathrm{IND}} - \theta^*\right)$ converges to a known distribution;

**Theorem: Indirect Estimator Consistency**

1. (Gourieroux et al., 1993) $\hat{\theta}_{\mathrm{IND}}$ is a consistent estimator of $\theta^*$, and $\sqrt{n}\left(\hat{\theta}_{\mathrm{IND}} - \theta^*\right)$ converges to a known distribution;

2. The parametric bootstrap CIs and HTs based on $\hat{\theta}_{\mathrm{IND}}$ are consistent.

**Theorem: Indirect Estimator Consistency**

1. (Gourieroux et al., 1993) $\hat{\theta}_{\mathrm{IND}}$ is a consistent estimator of $\theta^*$, and $\sqrt{n}\left(\hat{\theta}_{\mathrm{IND}} - \theta^*\right)$ converges to a known distribution;

2. The parametric bootstrap CIs and HTs based on $\hat{\theta}_{\mathrm{IND}}$ are consistent.

▶ Two levels of simulation in "PB+Indirect Estimator".

> **Theorem: Indirect Estimator Consistency**
>
> 1. (Gourieroux et al., 1993) $\hat{\theta}_{\text{IND}}$ is a consistent estimator of $\theta^*$, and $\sqrt{n}\left(\hat{\theta}_{\text{IND}} - \theta^*\right)$ converges to a known distribution;
> 2. The parametric bootstrap CIs and HTs based on $\hat{\theta}_{\text{IND}}$ are consistent.

▶ Two levels of simulation in "PB+Indirect Estimator".
  1. $s_i(\theta)$ used in $\hat{\theta}_{\text{IND}}$

> **Theorem: Indirect Estimator Consistency**
>
> 1. (Gourieroux et al., 1993) $\hat{\theta}_{\text{IND}}$ is a consistent estimator of $\theta^*$, and $\sqrt{n}\left(\hat{\theta}_{\text{IND}} - \theta^*\right)$ converges to a known distribution;
>
> 2. The parametric bootstrap CIs and HTs based on $\hat{\theta}_{\text{IND}}$ are consistent.

▶ Two levels of simulation in "PB+Indirect Estimator".
  1. $s_i(\theta)$ used in $\hat{\theta}_{\text{IND}}$
  2. $s_b(\hat{\theta}_{\text{IND}})$ used in PB

> **Theorem: Indirect Estimator Consistency**
>
> 1. (Gourieroux et al., 1993) $\hat{\theta}_{\mathrm{IND}}$ is a consistent estimator of $\theta^*$, and $\sqrt{n}\left(\hat{\theta}_{\mathrm{IND}} - \theta^*\right)$ converges to a known distribution;
> 2. The parametric bootstrap CIs and HTs based on $\hat{\theta}_{\mathrm{IND}}$ are consistent.

▶ Two levels of simulation in "PB+Indirect Estimator".
  1. $s_i(\theta)$ used in $\hat{\theta}_{\mathrm{IND}}$
  2. $s_b(\hat{\theta}_{\mathrm{IND}})$ used in PB
▶ The choice of $\Omega$ determines the asymptotic variance of $\hat{\theta}_{\mathrm{IND}}$.

**Definition (Adaptive indirect estimator)**

Let $\Omega = (S(\theta))^{-1}$. $S(\theta)$: sample covariance matrix of $\{s_i(\theta)\}_{i=1}^R$.
(Intuition: tolerate more difference in more uncertain directions.)

$$\hat{\theta}_{\mathrm{ADI}} := \underset{\theta \in \Theta}{\arg\min} \left\| s^* - \frac{1}{R} \sum_{i=1}^R s_i(\theta) \right\|_{(S(\theta))^{-1}}.$$

# (Novel) Adaptive Indirect Estimator

## Definition (Adaptive indirect estimator)

Let $\Omega = (S(\theta))^{-1}$. $S(\theta)$: sample covariance matrix of $\{s_i(\theta)\}_{i=1}^R$.
(Intuition: tolerate more difference in more uncertain directions.)

$$\hat{\theta}_{\mathrm{ADI}} := \underset{\theta \in \Theta}{\arg\min} \left\| s^* - \frac{1}{R} \sum_{i=1}^R s_i(\theta) \right\|_{(S(\theta))^{-1}}.$$

### Theorem: Consistency and asymptotic variance ($R \to \infty$)

1. $\hat{\theta}_{\mathrm{ADI}}$ is a consistent estimator of $\theta^*$, $\sqrt{n}\left(\hat{\theta}_{\mathrm{ADI}} - \theta^*\right)$ converges to a dist;

2. The parametric bootstrap CIs and HTs based on $\hat{\theta}_{\mathrm{ADI}}$ are consistent;

3. (Optimal asymptotic variance) For any well-behaved consistent estimator $\psi(s)$, we have $\mathrm{Var}\left(\lim_{n \to \infty} \sqrt{n}\left(\psi(s) - \theta^*\right)\right) \succeq \mathrm{Var}\left(\lim_{n \to \infty} \sqrt{n}\left(\hat{\theta}_{\mathrm{ADI}} - \theta^*\right)\right).$

# Inference of Parameters $(\mu^*, \sigma^*)$ of Normal

Figure: Sampling distributions of different estimates. Vertical line is median, '$*$' is true value.

# Inference of Parameters $(\mu^*, \sigma^*)$ of Normal

Figure: Sampling distributions of different estimates. Vertical line is median, '$*$' is true value.



| Method (95% CI) | Coverage | | Average width | |
|---|---|---|---|---|
| | $\mu^*$ | $\sigma^*$ | $\mu^*$ | $\sigma^*$ |
| PB (adaptive indirect) | 0.959 (0.006) | 0.951 (0.007) | 0.463 (0.003) | 0.580 (0.003) |
| PB (naïve percentile) | 0.697 (0.015) | 0.006 (0.002) | 0.311 (0.001) | 0.293 (0.001) |
| PB (simplified $t$) | 0.869 (0.011) | 0.817 (0.012) | 0.311 (0.001) | 0.293 (0.001) |
| PB (Ferrando et al., 2022) | 0.808 (0.012) | 0.371 (0.015) | 0.311 (0.001) | 0.293 (0.001) |
| PB (Efron's BC) | 0.854 (0.011) | 0.042 (0.006) | 0.298 (0.001) | 0.139 (0.002) |
| PB (automatic percentile) | 0.865 (0.011) | 0.126 (0.010) | 0.314 (0.001) | 0.261 (0.001) |
| Repro (Awan and Wang, 2023) | 0.989 (0.003) | 0.998 (0.001) | 0.599 (0.003) | 0.758 (0.005) |

$\mu^* = 1$, $\sigma^* = 1$, $L = 0$, $U = 3$, $R = 50$, $B = 200$, $(\sqrt{2})$-GDP.

▶ Compare rejection probabilities (level 0.05) to (Alabi and Vadhan, 2022) & Repro.



$\Delta = 2$, $\beta_0^* = -0.5$, $x_i \sim N(0.5, 1)$, $\varepsilon_i \sim N(0, 0.25)$, $R = 50$, $B = 200$, 1-GDP.

**Contributions:**

1. Prove consistency of PB (indirect estimator).
2. Propose an adaptive indirect estimator (ADI): consistent (PB), optimal asymp var.
3. Improve state-of-the-art DP PB (validity & efficiency).

---

**Limitations:**

1. Computationally expensive.

**Contributions:**

1. Prove consistency of PB (indirect estimator).
2. Propose an adaptive indirect estimator (ADI): consistent (PB), optimal asymp var.
3. Improve state-of-the-art DP PB (validity & efficiency).

---

**Limitations:**

1. Computationally expensive.
2. Requires regularity conditions (e.g., smoothness).

**Contributions:**

1. Prove consistency of PB (indirect estimator).
2. Propose an adaptive indirect estimator (ADI): consistent (PB), optimal asymp var.
3. Improve state-of-the-art DP PB (validity & efficiency).

---

**Limitations:**

1. Computationally expensive.
2. Requires regularity conditions (e.g., smoothness).
3. Additional techniques required for discrete settings (in the thesis).

**Contributions:**

1. Prove consistency of PB (indirect estimator).
2. Propose an adaptive indirect estimator (ADI): consistent (PB), optimal asymp var.
3. Improve state-of-the-art DP PB (validity & efficiency).

---

**Limitations:**

1. Computationally expensive.
2. Requires regularity conditions (e.g., smoothness).
3. Additional techniques required for discrete settings (in the thesis).

**Compare ADI to Repro: (both simulation-based)**

1. Repro is finite-sample valid with almost no assumptions while conservative.
2. ADI is an estimator, asymptotically optimal w/ more assumptions.
3. PB+ADI and Repro are state-of-the-art in different scenarios.

# Summary

- DP bootstrap: non-parametric; difficult DP analysis; but restricted in mechanisms;
- Repro: accepts all mechanisms; finite-sample valid; but conservative;
- PB+ADI: accepts all mechanisms; asymp valid & efficient; but needs smoothness.

Repro and PB+ADI are general-purpose methods that solve the clamping problem and outperform (Alabi and Vadhan, 2022) which only focused on linear regression.

|  | DP bootstrap | Repro | PB+ADI |
|---|---|---|---|
| **Data generating equation** | Not needed | Needed | Needed & Smooth |
| **DP mechanisms** | All (for DP guarantee); Additive-noise (for inference) | Easily sampled | Easily sampled & Smooth |
| **Inference** | Asymptotic; often conservative; requires a point estimator | Finite-sample; often conservative; no estimator | Asymptotic; efficient; provides an estimator |

**Future work**

▶ For Repro and PB+ADI, find an appropriate data generating equation?
  ▶ If there is none, consider non-parametric or semi-parametric models.
▶ Find the DP mechanism giving the optimal summary statistic $s$ for inference?
▶ For DP Bootstrap, we need more post-processing in addition to deconvolution if the original mechanism gives a biased estimator.

Daniel Alabi and Salil Vadhan. Hypothesis testing for differentially private linear regression. Advances in Neural Information Processing Systems, 35:14196–14209, 2022.

Jordan Awan and Zhanyu Wang. Simulation-based, finite-sample inference for privatized data. arXiv preprint arXiv:2303.05328, 2023.

Borja Balle, Gilles Barthe, and Marco Gaboardi. Privacy amplification by subsampling: Tight analyses via couplings and divergences. In Advances in Neural Information Processing Systems, volume 31, 2018.

Thomas Brawner and James Honaker. Bootstrap inference and differential privacy: Standard errors for free. Unpublished Manuscript, 2018.

Christian Covington, Xi He, James Honaker, and Gautam Kamath. Unbiased statistical estimation and valid confidence intervals under differential privacy. arXiv preprint arXiv:2110.14465, 2021.

Jinshuo Dong, Aaron Roth, and Weijie J. Su. Gaussian differential privacy. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 84(1):3–37, 2022.

Wenxin Du, Canyon Foot, Monica Moniot, Andrew Bray, and Adam Groce. Differentially private confidence intervals. arXiv preprint arXiv:2001.02285, 2020.

Bradley Efron. Empirical Bayes deconvolution estimates. Biometrika, 103(1):1–20, 2016.

Georgina Evans, Gary King, Margaret Schwenzfeier, and Abhradeep Thakurta. Statistically valid inferences from privacy-protected data. American Political Science Review, 117(4):12751290, 2023.

# References II

Cecilia Ferrando, Shufan Wang, and Daniel Sheldon. Parametric bootstrap for differentially private confidence intervals. In International Conference on Artificial Intelligence and Statistics, pages 1598–1618. PMLR, 2022.

Christian Gourieroux, Alain Monfort, and Eric Renault. Indirect inference. Journal of applied econometrics, 8 (S1):S85–S118, 1993.

Stéphane Guerrier, Elise Dupuis-Lozeron, Yanyuan Ma, and Maria-Pia Victoria-Feser. Simulation-based bias correction methods for complex models. Journal of the American Statistical Association, 114(525):146–157, 2019.

Wenxin Jiang and Bruce Turnbull. The Indirect Method: Inference Based on Intermediate Statistics – A Synthesis and Examples. Statistical Science, 19(2):239 – 263, 2004.

Nianqiao Ju, Jordan Awan, Ruobin Gong, and Vinayak Rao. Data augmentation MCMC for Bayesian inference from privatized data. In Advances in Neural Information Processing Systems, volume 36, 2022.

Antti Koskela, Joonas Jälkö, and Antti Honkela. Computing tight differential privacy guarantees using fft. In Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics, volume 108, pages 2560–2569. PMLR, 2020.

Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Algorithmic learning in a random world, volume 29. Springer, 2005.

Yue Wang, Daniel Kifer, Jaewoo Lee, and Vishesh Karwa. Statistical approximating distributions under differential privacy. Journal of Privacy and Confidentiality, 8(1), 2018.

Yue Wang, Daniel Kifer, and Jaewoo Lee. Differentially private confidence intervals for empirical risk minimization. Journal of Privacy and Confidentiality, 9(1), 2019.

Oliver Williams and Frank McSherry. Probabilistic inference and differential privacy. Advances in Neural Information Processing Systems, 23:2451–2459, 2010.

Min-ge Xie and Peng Wang. Repro samples method for finite-and large-sample inferences. arXiv preprint arXiv:2206.06421, 2022.

### Definition ($(\varepsilon, \delta)$-DP)

A mechanism $M : \mathcal{X}^n \to \mathcal{Y}$ is $(\varepsilon, \delta)$-DP if for any neighboring datasets $D \simeq D' \in \mathcal{X}^n$, and any measurable set $S \subseteq \mathcal{Y}$, the following inequality holds:

$$\Pr[M(D) \in S] \leq e^{\varepsilon} \Pr[M(D') \in S] + \delta.$$

### Definition (tradeoff function & $f$-DP)

Consider the hypothesis test $H_0 : Y \sim P$ versus $H_1 : Y \sim Q$. For any rejection rule $\phi(Y)$, $\alpha_\phi$ is the type I error and $\beta_\phi$ is the type II error. The tradeoff function is

$$T_{P,Q}(\alpha) := \inf_\phi \{ \beta_\phi \mid \alpha_\phi \leq \alpha \}.$$

$M$ is $f$-DP if $T_{M(D),M(D')}(\alpha) \geq f(\alpha)$ for any $\alpha$ and datasets $D, D'$ with $D \simeq D'$.

**Primal-dual conversion** $f$-DP $\Leftrightarrow$ $(\varepsilon, \delta)$-DP $\forall \varepsilon \geq 0$ with $\delta(\varepsilon) = 1 + f^*(-e^{\varepsilon})$.
▶ $f_{\varepsilon,\delta}(\alpha) := \max\{0, \ 1 - \delta - e^{\varepsilon}\alpha, \ e^{-\varepsilon}(1 - \delta - \alpha)\}$-DP is equivalent to $(\varepsilon, \delta)$-DP.

# Privacy Guarantees for Mixture Mechanism

▶ If $M$ randomly releases the output of $M_i$ with probability $p_i$, and $M_i$ is $f_i$-DP for $i \in [k]$, then $M$ is $f_{\mathrm{mix}}$-DP. $f_{\mathrm{mix}} = \left( \sum_{i=1}^{k} \left( p_i f_i \circ (f_i')^{-1} \right) \right) \circ \left( \sum_{i=1}^{k} p_i (f_i')^{-1} \right)^{-1}$

▶ $\underline{p} = (p_1, p_2, p_3) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. $f_1, f_2, f_3$ correspond to 1-GDP, 2-GDP, and 3-GDP.

# Privacy Analysis of DP Bootstrap: An Example

▶ Consider the Gaussian mechanism $M$ with 1-GDP (dashed curve). $M \circ \text{boot}$ satisfies $f_{\text{boot}}$-DP (solid opaque curve). The transparent curves are for testing $M(D)$ vs $M(D')$ where $D = (a, 0, \ldots, 0)$, $D' = (a - 1, 0, \ldots, 0)$, $M(D) = \frac{1}{n} \sum_{i=1}^{n} x_i + \xi$, $D = (x_1, x_2, \ldots, x_n)$, $\xi \sim \mathcal{N}(0, \frac{1}{n^2})$.

▶ The dashed and dotted dashed lines are misused as lower bounds in Brawner and Honaker (2018) and Koskela et al. (2020).

# Accuracy of DP Bootstrap with Gaussian Mechanism

| Unbiased estimate of $\theta$ | Variance of the estimate |
|---|---|
| *Sample mean (non-private):* $\hat{\theta}_1 = \bar{X}$ | $\mathrm{Var}(\hat{\theta}_1) = \frac{\sigma_x^2}{n}$ |
| *Sample mean (Gaussian mechanism):* $\hat{\theta}_2 = \bar{X} + \xi$ where $\xi \sim \mathcal{N}(0, \frac{1}{\mu^2 n^2})$ | $\mathrm{Var}(\hat{\theta}_2) = \frac{\sigma_x^2}{n} + \frac{1}{\mu^2 n^2}$ |
| *Bootstrap (non-private):* $\hat{\theta}_3 = \bar{X}$ | $\mathrm{Var}(\hat{\theta}_3) = \frac{\sigma_x^2}{n}$ |
| *DP bootstrap (Gaussian mechanism):* $\hat{\theta}_4 = \tilde{X}'$ where $\xi_b \sim \mathcal{N}(0, \frac{(2-2/e)B}{\mu^2 n^2})$ | $\mathrm{Var}(\hat{\theta}_4) = \frac{1+1/B-1/(nB)}{n}\sigma_x^2 + \frac{(2-2/e)}{\mu^2 n^2}$ |

| Unbiased estimate of $\mathrm{Var}(\hat{\theta}_i)$ | Variance of the estimate |
|---|---|
| $\widehat{\mathrm{Var}}(\hat{\theta}_1) = \frac{s_X^2}{n}$ | $\mathrm{Var}(\widehat{\mathrm{Var}}(\hat{\theta}_1)) \in O(\frac{1}{n^3})$ |
| $\widehat{\mathrm{Var}}(\hat{\theta}_2) = \frac{s_X^2 + \xi}{n} + \frac{1}{\mu^2 n^2}$ where $\xi \sim \mathcal{N}(0, \frac{1}{\mu^2 n^2})$ | $\mathrm{Var}(\widehat{\mathrm{Var}}(\hat{\theta}_2)) \in O(\frac{1}{n^3} + \frac{1}{\mu^2 n^4})$ |
| $\widehat{\mathrm{Var}}(\hat{\theta}_3) = \frac{n}{n-1}\tilde{s}_B^2$ | $\mathrm{Var}(\widehat{\mathrm{Var}}(\hat{\theta}_3)) \in O(\frac{1}{n^2 B} + \frac{1}{n^3})$ |
| $\widehat{\mathrm{Var}}(\hat{\theta}_4) = \frac{nB+n-1}{B(n-1)}\tilde{s}_B^2 - \frac{(2-2/e)B}{n(n-1)\mu^2}$ | $\mathrm{Var}(\widehat{\mathrm{Var}}(\hat{\theta}_4)) \in O(\frac{1}{n^2 B} + \frac{1}{n^3 \mu^2} + \frac{B}{n^4 \mu^4} + \frac{1}{n^3})$ |

# DP Bootstrap: Deconvolution

▶ We choose to use `deconvolveR` (Efron, 2016) based on Empirical Bayes since it performs the best in our settings without tuning its hyper-parameters.

▶ For the model $Y = X + e$, `deconvolveR` assumes that $Y$ and $X$ are distributed discretely with the sizes of their supports $|\mathcal{Y}| = k$ and $|\mathcal{X}| = m$.

▶ It models the distribution of $X$ by $f(\alpha) = e^{Q\alpha}/c(\alpha)$ where $Q$ is an $m \times p$ structure matrix with values from the natural spline basis with order $p$, $ns(\mathcal{X}, p)$, and $\alpha$ is the unknown $p$-dimensional parameter vector; $c(\alpha)$ is the divisor necessary to make $f$ sum to 1.

▶ The estimation of the distribution of $X$ is obtained through the estimation of $\alpha$: It estimates $\alpha$ by maximizing a penalized log-likelihood $m(\alpha) = l(Y; \alpha) - s(\alpha)$ with respect to $\alpha$ where $s(\alpha)$ is the penalty term, and $l(Y; \alpha)$ is the log-likelihood function of $Y$ derived from $f(\alpha)$ and the known distribution of $e$.
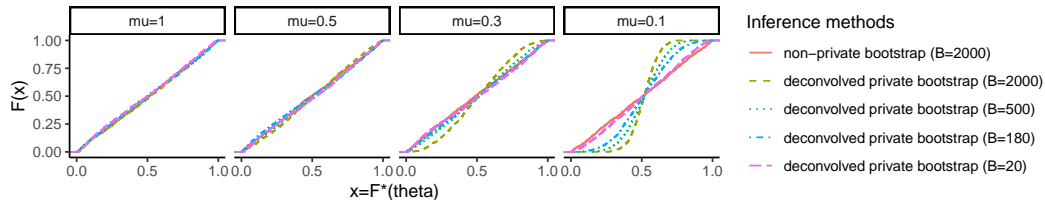
# DP CIs for Normal: Compare DP Bootstrap to NoisyVar (Du et al., 2020)

Table: Coverage and width of CIs with different privacy guarantees. Confidence level is 90%. The standard error estimated from 2000 replicates is in parenthesis. $B \propto n\mu^2$.

| Privacy | Method | Coverage | CI width |
|---------|--------|----------|----------|
| N/A | Bootstrap (B=2000) | 0.905 (7e-3) | 0.014 (6e-6) |
| 1-GDP | DP bootstrap (B=2000) | 0.896 (7e-3) | 0.014 (1e-5) |
| | NoisyVar | 0.803 (9e-3) | 0.011 (7e-6) |
| 0.5-GDP | DP bootstrap (B=500) | 0.898 (7e-3) | 0.014 (2e-5) |
| | NoisyVar | 0.806 (9e-3) | 0.011 (7e-6) |
| 0.3-GDP | DP bootstrap (B=180) | 0.901 (7e-3) | 0.015 (3e-5) |
| | NoisyVar | 0.804 (9e-3) | 0.011 (7e-6) |
| 0.1-GDP | DP bootstrap (B=20) | 0.962 (4e-3) | 0.020 (1e-4) |
| | NoisyVar | 0.819 (9e-3) | 0.012 (7e-6) |

▶ Coverage check for all confidence levels.

# Repro: Over-coverage

Table: Relative width due to over-coverage for the normal mean with known variance, when the nominal level is $1 - \alpha = 0.95$, and the over-coverage level is $1 - \alpha^* = (1 - \alpha)^{1/d}$.

| Dimension $d$ | 1 | 2 | 5 | 10 | 100 | 1000 |
|---|---|---|---|---|---|---|
| Relative width | 1 | 1.14 | 1.31 | 1.43 | 1.77 | 2.07 |

Table: 95% confidence intervals for private Bernoullis with unknown $n$. The first row uses Mahalanobis depth, and the second row uses an approximate pivot. For both intervals, an initial $(1 - 10^{-4})$-CI for $n$ is used to reduce the nuisance parameter search. Parameters for the simulation are $n^* = 100$, $p^* = 0.2$, $\varepsilon = 1$, $R = 200$.

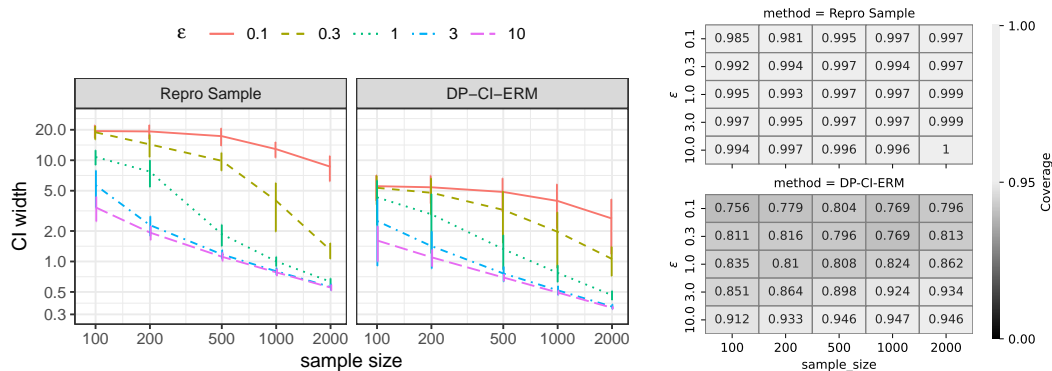| | Coverage | Width |
|---|---|---|
| Mahalanobis Depth | 0.980 (0.004) | 0.197 (0.001) |
| Approximate Pivot | 0.949 (0.007) | 0.163 (0.001) |

Figure: Width and coverage for the confidence intervals of $\beta_1$ in logistic regression with repro and DP-CI-ERM Wang et al. (2019). Parameters for this simulation are $a^* = b^* = 0.5$, $\beta_0^* = 0.5$, $\beta_1^* = 2$, $R = 200$, $\alpha = 0.05$, and the results were averaged over 1000 replicates.

$$\hat{\theta}_{\mathrm{DP}}(D; u) = \arg\min_{\theta \in \Theta} \left( \hat{\mathcal{L}}(\theta; D) + \frac{\gamma}{2n} \theta^{\mathsf{T}} \theta + \frac{u^{\mathsf{T}} \theta}{n} \right), \quad f(u; \varepsilon, \Delta) \propto \exp\left( -\frac{\varepsilon q}{\Delta} \|u\|_\infty \right).$$

# Logistic Regression: Different Test Statistics in Repro

Table: Average width for the confidence intervals of $\beta_1$ in logistic regression using repro with the Mahalanobis depth on different summary statistics $s$. $T_{\text{pivot}} :=$

$$\sqrt{n}\left(\hat{l}(\beta^*; D_{\theta^*}) + \text{Cov}(V)\right)^{-\frac{1}{2}}\left(\left(H^* + \frac{\gamma}{n}\right)\hat{\theta}_{\text{DP}} - H^*\beta^*\right), \ H^* := \frac{1}{2}(\hat{l}(\hat{\theta}_{\text{DP}}; D_{\hat{\theta}^*_{\text{DP}}}) + \hat{l}(\beta^*; D_{\theta^*})).$$

| $s = (\hat{\theta}_{\text{DP}}, \tilde{z}, \widetilde{z^2})$ | $n = 100$ | $n = 200$ | $n = 500$ | $n = 1000$ | $n = 2000$ |
|---|---|---|---|---|---|
| $\varepsilon = 0.1$ | 19.430 (0.089) | 19.252 (0.099) | 17.306 (0.109) | 12.870 (0.072) | 8.622 (0.077) |
| $\varepsilon = 0.3$ | 18.877 (0.091) | 14.335 (0.114) | 9.878 (0.064) | 3.975 (0.064) | 1.291 (0.007) |
| $\varepsilon = 1$ | 10.762 (0.057) | 7.727 (0.073) | 1.862 (0.014) | 1.003 (0.004) | 0.630 (0.002) |
| $\varepsilon = 3$ | 5.678 (0.071) | 2.287 (0.016) | 1.176 (0.004) | 0.801 (0.002) | 0.560 (0.001) |
| $\varepsilon = 10$ | 3.426 (0.030) | 1.931 (0.010) | 1.115 (0.004) | 0.781 (0.002) | 0.553 (0.001) |

| $s = (T_{\text{pivot}}, \tilde{z}, \widetilde{z^2})$ | $n = 100$ | $n = 200$ | $n = 500$ | $n = 1000$ | $n = 2000$ |
|---|---|---|---|---|---|
| $\varepsilon = 0.1$ | 19.240 (0.108) | 19.337 (0.097) | 18.487 (0.105) | 14.789 (0.106) | 7.000 (0.097) |
| $\varepsilon = 0.3$ | 19.243 (0.087) | 15.533 (0.134) | 8.234 (0.091) | 2.577 (0.032) | 1.148 (0.005) |
| $\varepsilon = 1$ | 9.939 (0.084) | 4.613 (0.062) | 1.594 (0.008) | 0.955 (0.003) | 0.617 (0.002) |
| $\varepsilon = 3$ | 3.309 (0.033) | 1.905 (0.009) | 1.118 (0.003) | 0.782 (0.002) | 0.553 (0.001) |
| $\varepsilon = 10$ | 2.381 (0.012) | 1.665 (0.006) | 1.058 (0.003) | 0.762 (0.002) | 0.545 (0.001) |

# Indirect Estimator (Gourieroux, Monfort, and Renault, 1993)

- Private statistics:

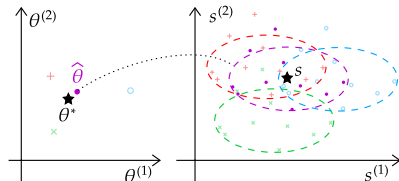$$s := \arg\min_{\beta} \rho(\beta; D, u_{\mathrm{DP}}), \ D := G(\theta^*; u).$$

  - DP mechanism $\rho$ contains extra uncertainty $u_{\mathrm{DP}} \sim F_{\mathrm{DP}}$,
  - e.g., Gaussian Mechanism, Objective perturbation.
- Fix the randomness $(u^r, u_{\mathrm{DP}}^r)$ in generating $D^r(\theta) := G(\theta, u^r)$ and

$$s^r(\theta) := \arg\min_{\beta} \rho(\beta; D^r(\theta), u_{\mathrm{DP}}^r).$$

- Find the $\theta$ generating $s^r(\theta)$ most similar to $s$.

Definition (Indirect estimator)

$$\hat{\theta}_{\mathrm{IND}} := \arg\min_{\theta \in \Theta} \left\| s - \frac{1}{R} \sum_{r=1}^{R} s^r(\theta) \right\|_{\Omega_n}.$$

# Indirect Estimator: Asymptotic Distributions

- $\rho_n(\beta; D, u_{\mathrm{DP}}) \xrightarrow{\mathrm{P}} \rho_\infty(\beta; F_u, F_{\mathrm{DP}}, \theta^*),$

- $b(\theta) := \mathrm{argmax}_{\beta \in \mathbb{B}} \rho_\infty(\beta; F_u, F_{\mathrm{DP}}, \theta), \ \beta^* := b(\theta^*),$

- $B^* := \frac{\partial b(\theta^*)}{\partial \theta}, \ J^* := -\frac{\partial^2 \rho_\infty(\beta^*; F_u, F_{\mathrm{DP}}, \theta^*)}{(\partial \beta)(\partial \beta^\intercal)},$

- $\sqrt{n}\left(\frac{\partial \rho_n(\beta^*; D, u_{\mathrm{DP}})}{\partial \beta}\right) \xrightarrow{\mathrm{d}} F^*_{\rho, u, \mathrm{DP}}, \quad \Omega_n \to \Omega.$

Let $v_i \overset{\mathrm{iid}}{\sim} F^*_{\rho, u, \mathrm{DP}}$, $\Sigma^* := \mathrm{Var}[(J^*)^{-1} v_0] = \mathrm{Var}\left(\lim_{n \to \infty} \sqrt{n}(s - b(\theta^*))\right)$, $\Omega^* := (\Sigma^*)^{-1}$.

$$\sqrt{n}(\hat{\theta}_{\mathrm{IND}} - \theta^*) \xrightarrow{d} ((B^*)^\intercal \Omega B^*)^{-1} (B^*)^\intercal \Omega (J^*)^{-1} \left(v_0 - \frac{1}{R} \sum_{i=1}^R v_i\right).$$

$$\sqrt{n}(\hat{\theta}_{\mathrm{ADI}} - \theta^*) \xrightarrow{d} ((B^*)^\intercal \Omega^* B^*)^{-1} (B^*)^\intercal \Omega^* (J^*)^{-1}(v_0 - \mathbb{E}(v_0)).$$

$$\mathrm{Var}\left(\lim_{n \to \infty} \sqrt{n}(\hat{\theta}_{\mathrm{IND}} - \theta^*)\right) \succeq \mathrm{Var}\left(\lim_{n \to \infty} \sqrt{n}(\hat{\theta}_{\mathrm{ADI}} - \theta^*)\right) = \textcolor{red}{\left((B^*)^\intercal (\Sigma^*)^{-1} B^*\right)^{-1}}.$$

- Test statistic $\hat{\tau}$. Auxiliary scale of test statistic $\hat{\sigma}$.
- Let $\hat{\hat{\xi}}_{(j)}$ be the $j$th order statistic of $\left\{ \frac{\hat{\tau}(s_b) - \tau(\hat{\theta})}{\hat{\sigma}(s_b)} \right\}_{b=1}^{B}$.
- CI for $\tau(\theta^*)$ is $\left[ \hat{\tau}(s) + \hat{\hat{\xi}}_{(\lfloor (B+1)\alpha/2 \rfloor)}\hat{\sigma}(s), \ \hat{\tau}(s) + \hat{\hat{\xi}}_{(1+B-\lfloor (B+1)\alpha/2 \rfloor)}\hat{\sigma}(s) \right]$.

We want to choose $\hat{\tau}$ and $\hat{\sigma}$ such that $\frac{\hat{\tau}(s_b) - \tau(\hat{\theta})}{\hat{\sigma}(s_b)}$ has mean 0 and variance 1.

- Note that $b(\theta^*) = \lim_{n \to \infty} s^*$ and $\Sigma(\theta^*) = \mathrm{Var}\left( \lim_{n \to \infty} \sqrt{n}(s - b(\theta^*)) \right)$.
- Let $\hat{\theta} := \hat{\theta}_{\mathrm{ADI}}$, $\hat{\theta}_b := \hat{\theta}(s_b)$. Set the test statistic as $\hat{\tau}(s_b) := \eta_1(\hat{\theta}(s_b))$. We use $\hat{\sigma}(s_b)$ to estimate the asymptotic standard deviation of $\hat{\tau}(s_b)$, where

$$\hat{\sigma}(s_b) := \frac{1}{\sqrt{n}} \left( \frac{\partial \eta_1}{\partial \theta}(\hat{\theta}_b) \left( \left( \frac{\partial b}{\partial \theta}(\hat{\theta}_b) \right)^{\mathsf{T}} \Sigma(\hat{\theta}_b)^{-1} \frac{\partial b}{\partial \theta}(\hat{\theta}_b) \right)^{-1} \left( \frac{\partial \eta_1}{\partial \theta}(\hat{\theta}_b) \right)^{\mathsf{T}} \right)^{\frac{1}{2}}.$$

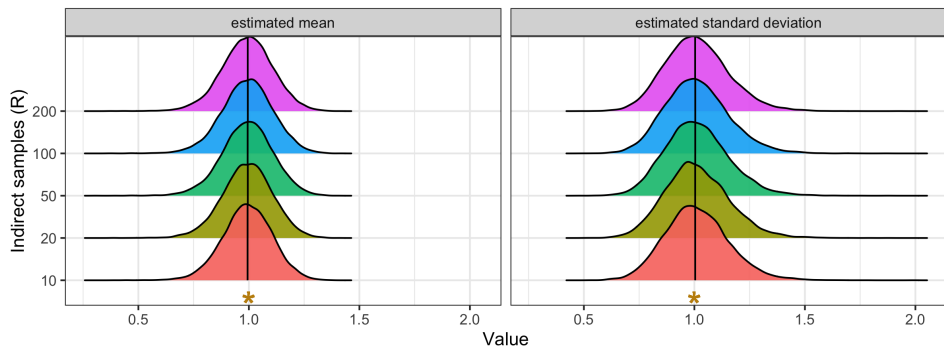# Indirect Estimator with Different $R$



Figure: Comparison of the sampling distribution of the adaptive indirect estimates $\hat{\theta}_{\mathrm{ADI}}$ under different settings of the number of generated samples $R = 10, 20, 50, 100, 200$ in the normal distribution setting.
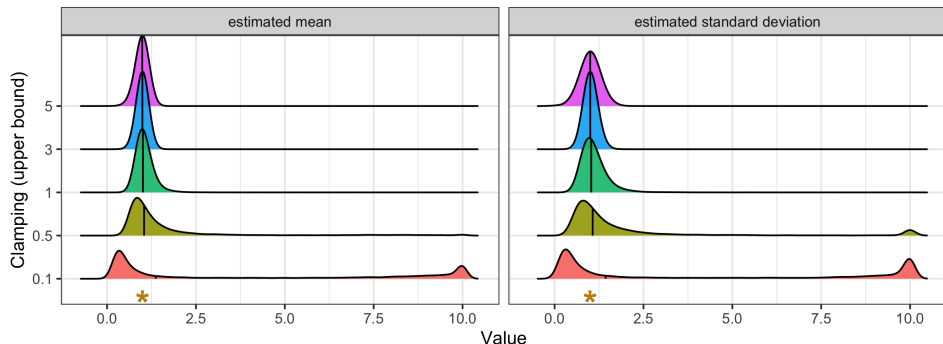
Figure: Comparison of the sampling distribution of the adaptive indirect estimates $\hat{\theta}_{\mathrm{ADI}}$ under different settings of the clamping parameter $U = 0.1, 0.5, 1, 3, 5$ in the normal distribution setting.