

# Flexible Formal Privacy for Public Data Curation

Jeremy Seeman

(based on joint work with A. Slavkovic, M. Reimherr, and others)

Michigan Institute for Data Science (MIDAS) and Institute for Social Research (ISR)

University of Michigan

Contact: [jhseeman@umich.edu](mailto:jhseeman@umich.edu)

November 1, 2023

## Goal

Discuss how flexible formal privacy methods can help address privacy concerns while limiting new operational barriers to social science data collection.

In this talk, I'll discuss...

- 1 What formal privacy (FP) methods, such as differential privacy (DP) offers for public data curation.
- 2 Why applying DP to survey data is especially challenging.
- 3 What strategic relaxations of DP offer FP guarantees without interfering with survey operations.
- 4 How to design FP mechanisms which accommodate reproducible inferences from social science data.

Large-scale data and computing makes reconstructing personal information from published datasets easier than ever before.

- Traditional work: statistical disclosure limitation (SDL) techniques evaluate *properties of datasets* to determine individuals' risk of reidentification [HDF+2012]
- Modern problem: *any* statistic can be used to help reidentify individuals! [DN2003,K2009]

2010 Census data is highly susceptible to reconstruction [A2021]:

- Tabular summaries from the U.S. Census can be used to reconstruct potential populations which conform to these summaries
- Confirmed reidentifications for  $\approx 52$  million respondents,  $\approx 16\%$  population

Record Linkage Summary from Commercial and CEF Record Sources				
PIK, Block, Age, Sex Record Linkage <b>Source</b>	Available Records	Records with PIK, Block, Sex, and Age	Putative Re-identifications using <b>Source</b>	Confirmed Re-identifications
Commercial	413,137,184	286,671,152	137,709,807	52,038,366

Figure: Reconstruction attack results from third-party data (source: U.S. Census)

Reconstruction attacks have tangible harms:

- Ex1: Reidentifying trans youth [KF22]
- Ex2: FERPA violations [C22]

As a result, data curators now turn to privacy-enhancing technologies (PETs):

- U.S. Census Bureau's disclosure avoidance modernization efforts satisfy Title 13 and 26 requirements.
- Biden's executive order (last Monday) proposes broad adoption of PETs for statistical agencies.



Source: <https://www.thenation.com/article/activism/texas-gender-affirming-care/>



Source: <https://educlove.com/best-free-edx-courses-with-certificates/>

WH.GOV



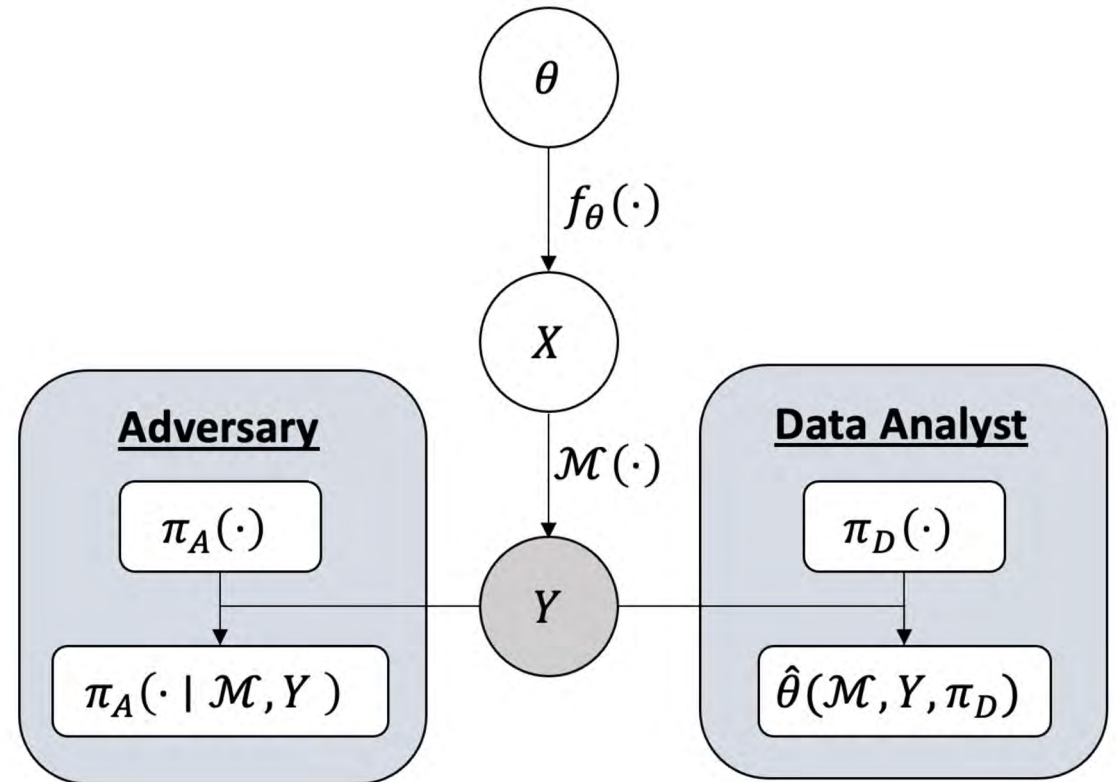
OCTOBER 30, 2023

FACT SHEET: President Biden Issues Executive Order on Safe, Secure, and Trustworthy Artificial Intelligence



## Notation:

- $\theta \in \Theta$ : parameter
- $f_\theta(\cdot)$ : model
- $X \in \mathcal{X}$ : database
- $Y \in \mathcal{Y}$ : mechanism output
- $\mathcal{M}$ : randomized algorithm
- $\pi_A(\cdot)$ : adversary prior
- $\pi_D(\cdot)$ : data analyst prior



## Privacy Goal

Enable inferences about global parameters  $\theta$  without leaking information about confidential data  $X$ .

Formal privacy makes quantifying disclosure risks *inherent to the mechanism*  $\mathcal{M}$ .

Desiderata:

- **Methodological transparency:** knowledge of a release strategy should not disclose additional confidential information.
- **Robustness to post-processing:** additional data processing operations shouldn't degrade privacy guarantees.
- **Privacy accounting:** privacy preservation should be quantified and accounted for in different statistical tasks.

DP [DMNS2006] is a popular framework for releasing statistical results with relative robustness to individuals' data contributions:

- $\mathcal{X} \triangleq$  sample space of 1 individual's data
- $(\mathcal{Y}, \mathcal{F}) \triangleq$  output space
- $\mathcal{M} \triangleq \{\mu_X \mid X \in \mathcal{X}^n\}$  release mechanism

### Differential privacy [DMNS2006]

A mechanism  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -DP if, for all  $B \in \mathcal{F}$  and adjacent  $X, X' \in \mathcal{X}^n$  (i.e.,  $X, X'$  differing on one record):

$$\mu_X(B) \leq e^\epsilon \mu_{X'}(B) + \delta$$

When  $\delta = 0$ , we say  $\mathcal{M}$  satisfies  $\epsilon$ -DP. Notes:

- DP mechanisms require randomized noise with magnitude inversely proportional to  $\epsilon$ .
- Privacy utility trade-off: smaller  $\epsilon$ , greater privacy, less data utility.



## DP changes the unit of analysis for disclosure avoidance.

- SDL: *absolute* disclosure risks as a property of published statistics.
- DP: *relative* disclosure risks as a property of a schema ( $\mathcal{X}$ ) and mechanism.

Testing interpretations: suppose, WLOG, we want to identify the first record  $X_1$ :

$$H_0 : X_1 = v_0, \quad H_1 : X_1 = v_1$$

- Any level- $\alpha$   $(\epsilon, \delta)$ -DP has Type II error *at least* [WZ2010]:

$$f_{\epsilon, \delta}(\alpha) \triangleq \max \{0, 1 - \delta - \alpha e^\epsilon, e^{-\epsilon}(1 - \delta - \alpha)\}.$$

- When  $\delta = 0$ , Bayes factors are bounded within  $[e^{-\epsilon}, e^\epsilon]$ . [KM2012]

Private selection:

- Goal: minimize loss function  $L_X$  while satisfying  $\epsilon$ -DP

$$L_X : \mathcal{X}^n \times \mathcal{Y} \mapsto [0, \infty]$$

- Key ingredient: bounded sensitivity of  $L_X$ . For all adjacent  $X, X' \in \mathcal{X}^n$ :

$$|L_X(y) - L_{X'}(y)| \leq \Delta_L < \infty.$$

### Exponential mechanism [MT2007]

A sample from density  $f_X$  with the form:

$$f_X(y) \propto \exp\left(-\frac{\epsilon L_X(y)}{2\Delta_L}\right),$$

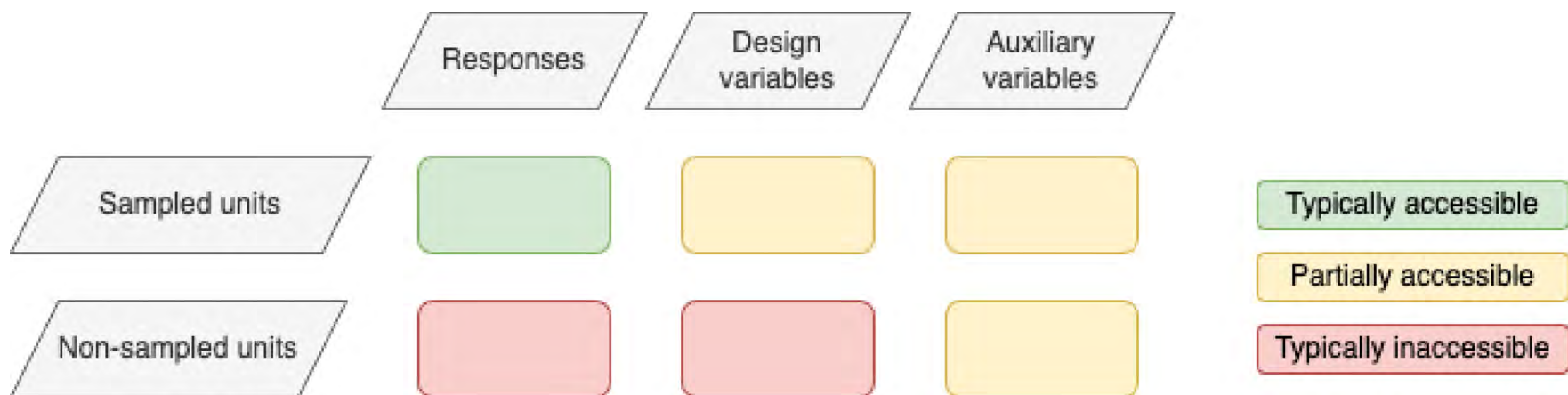
with respect to a common base measure  $\nu(y)$  over  $(\mathcal{Y}, \mathcal{F})$  satisfies  $\epsilon$ -DP.

## Open issues

DP requires additional randomness for every statistic derived from the confidential data, **but this doesn't happen in practice!**

- DP *only protects against information leakage attributable to the mechanism form* [KS2008,KS2014]
- Most DP analyses *intentionally ignore the joint effects of public information released about the same confidential dataset.*
- DP treats all public information as *equally disclosive, which may not be practically relevant.*

- “Adjacent databases” definitions tend to be over-inclusive for surveys.
  - Requires protecting auxiliary data that researchers may not be able to access.
  - Obstructs “secrecy of the sample” based on pathological worst-case scenarios.
- Pure DP analyses prevent deterministic survey design decision-making.
  - Making survey methodology public can violate DP [SB21,BDG+22]
  - Data-dependent sampling designs need to be randomized to ensure DP, which can be economically and/or logistically infeasible [BDG+22,D23].
  - Ex: “With non-zero probability, DP data-dependent sampling rates can take any value between 0 and 1.”



Decisions surrounding how to use all PETs, not just DP, are contentious:

- Social science researchers tend to strongly endorse or strongly oppose DP because it intentionally degrades data quality in service of a social value...
- ...making holistic policy decisions about how to implement *any* PETs, not just DP, particularly challenging [S23]<sup>1</sup>.

### Goal: strategic relaxations of DP for survey methodology

Provide FP methodology that simultaneously...

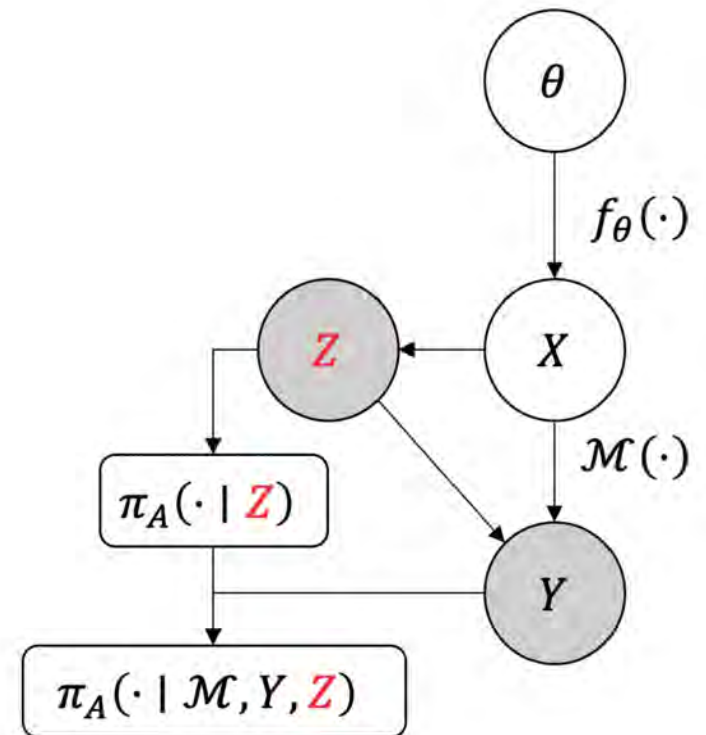
- 1 Enables formal privacy guarantees which protect sensitive information...
- 2 ...while minimizing interference with survey operations.

---

<sup>1</sup>Seeman, 2023. “Better Privacy Theorists for Better Data Stewards.” (*forthcoming*) *JPC*.

We define *public information* as any information dependent on  $X$  without privacy-preserving noise (call this r.v.  $Z$ ).

The relationship between  $X$  and  $Z$  may be deterministic or probabilistic.



Public information emerges from external societal forces:

- At implementation time: laws, social norms, and other established contextual information flows, decision-making process details.
- Before implementation: past releases and time dependence, establishment of public domains / common knowledge.

Examples:

- Summary statistics (ex: tabular summaries [GM2020])
- Database structure (ex: manifolds [RBS2021])
- Autocorrelation structure (ex: spatiotemporal models [Q2020])
- Known phenomenological structure (ex: genomic data [AAU2020])
- Sampling methodology (ex: surveys [SB2021])
- Fitness-for-use statistics (ex: linear query errors [XDW+2021])

## General approach:

- 1 Release essential statistics as-is when contextually necessary.
- 2 Adjust sensitive, granular statistics with formal privacy protections given essential statistics.

## Technical change:

- Standard DP: marginal distributions of data releases given adjacent databases should be close.
- **Our proposal:** conditional distributions given “adjacent databases” *and* public information should be close.
  - NB: adjacent databases now means conditional information about individual records,  $s_1, s_2 \in \mathbb{S}_{\text{pairs,DP}}$ , making the new testing problem

$$H_0 : s_1, \quad H_1 : s_2$$

- Example:  $s_1 =$  “the first record in  $X$  has value  $v_1 \in \mathcal{X}$ ” ( $s_2$  similar with  $v_2$ ).



We propose DP-style privacy guarantees as properties on conditional distributions  $X | Z$ , extending [KM2014]:

Def:  $\epsilon$ -TP [SSR2022] <sup>1</sup>

For each  $z \in \mathcal{Z}$ , let  $\mathbb{D}_{\text{DP}_z}$  be a collection of conditional distributions for  $X | Z = z$  indexed by  $\theta_z \in \Theta_Z$ . We say  $Y$  satisfies  $\epsilon$ -TP if for all  $z \in \mathcal{Z}$  and  $B \in \mathcal{F}_Y$ , for all distributions  $\theta_z \in \mathbb{D}_{\text{DP}_z}$ , and for all  $(s_1, s_2) \in \mathbb{S}_{\text{pairs,DP}_z}$ , where:

$$\mathbb{S}_{\text{pairs,DP}_z} \triangleq \{(s_1, s_2) \in \mathbb{S}_{\text{pairs,DP}} \mid \mathbb{P}(s_i | \theta_z) \notin \{0, 1\} \quad \forall i \in \{1, 2\}, \theta_z \in \mathbb{D}_{\text{DP}_z}\},$$

we have:

$$\begin{cases} \mathbb{P}(Y \in B | s_1, \theta_z) \leq e^\epsilon \mathbb{P}(Y \in B | s_2, \theta_z) \\ \mathbb{P}(Y \in B | s_2, \theta_z) \leq e^\epsilon \mathbb{P}(Y \in B | s_1, \theta_z). \end{cases}$$

<sup>1</sup>Seeman et al, 2022. "Formal Privacy for Partially Private Data." *Under Revision at JMLR*.

Suppose we observe  $X_1, \dots, X_n \in [-\Delta/2, \Delta/2]$ , and we consider data generating processes of the form for estimating  $\bar{X}$ :

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2), \quad \mu \in [-\Delta/2, \Delta/2], \quad \sigma \in (0, \infty)$$

Suppose we assume  $Z$  is jointly MVN with  $\bar{X}$  by introducing new covariance parameters:

$$\begin{pmatrix} \bar{X} \\ X_1 \\ Z \end{pmatrix} \sim N \left( \begin{pmatrix} \mu \\ \mu \\ \mu_Z \end{pmatrix}, \begin{pmatrix} \frac{\sigma^2}{n} & \Sigma_{XVZ}^T \\ \Sigma_{XVZ} & \Sigma_{VZ} \end{pmatrix} \right)$$

- $\Sigma_{XVZ}$  and  $\Sigma_{VZ}$  determines the strength of relationships between  $Z$  (public information),  $\bar{X}$  (statistic of interest) and  $X_1$  (one record).
- If  $\Sigma_{XVZ}$  isn't full rank, then  $Z$  is a linear combination of  $X_1, \dots, X_n \implies$  the support of  $X_1, \dots, X_n$  has smaller dimension!

Suppose we have  $j \in J$  strata and perform SRS without replacement within each strata. Existing approaches either...

- violate DP (like Neyman allocation) [BDG+22]...
- ... or use public strata sizes and ignore joint privacy concerns [LBG+23].

Example application of  $\epsilon$ -TP:

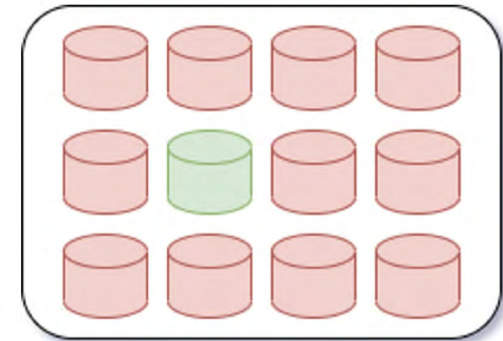
- Public information: population and sample strata sizes,  $z \triangleq \{(N_j, n_j)\}_{j=1}^J$
- Data generating distributions:  $\theta_z \in \Theta_z$  indexes all SRS without replacement distributions given population and sample strata sizes  $z$ .
- Protected attributes: within-strata responses but not strata inclusion (example interpretation below).

$$\forall j \in [J], i \in [n_j], \frac{\mathbb{P}(\text{Output} \mid \text{Unit } i \text{ in strata } j \text{ has response } y_1)}{\mathbb{P}(\text{Output} \mid \text{Unit } i \text{ in strata } j \text{ has response } y_2)} \leq \exp(\epsilon)$$

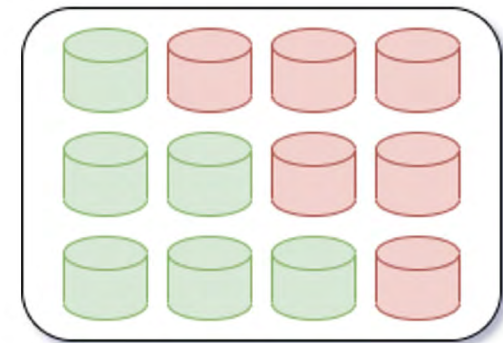
$\epsilon$ -TP interpolates between the units of analysis for statistical disclosure limitation (SDL) and DP.

- The bad news:
  - Guarantees non-uniform across schema  $\implies$  certain database reconstructions will be possible
  - Limits generalizability of composition results
- The good news:
  - For the partial schema,  $\epsilon$ -TP maintains similar desirable properties to  $\epsilon$ -DP
  - SDL and  $\epsilon$ -TP could be satisfied simultaneously

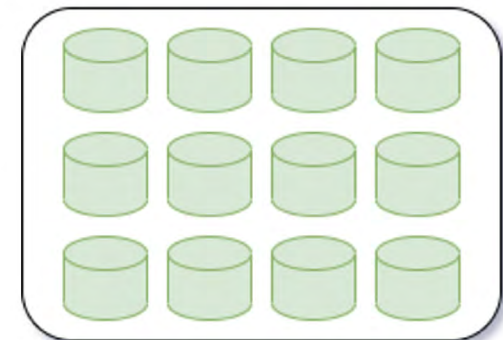
SDL  
(equivalent  
databases)



$\epsilon$ -TP  
(partial  
schema)



$\epsilon$ -DP  
(complete  
schema)



Two key differences when using  $\epsilon$ -TP versus  $\epsilon$ -DP:

- **Sensitivity inflation:**  $\epsilon$ -TP requires more noise than  $\epsilon$ -DP as...
  - Released statistics depend more on  $Z$  than  $Y$
  - Space of possible data generating scenarios  $\Theta$  grows.
- **Prior regularization:** use  $Z = z$  in the limit as  $\epsilon \rightarrow 0$ .
  - Standard DP: as  $\epsilon \rightarrow 0$ , release uniform noise
  - $\epsilon$ -TP: as  $\epsilon \rightarrow 0$ , release  $\theta \mid Z = z$ .

We can find approximate solutions to optimization problems with loss functions  $L_X$  while satisfying  $\epsilon$ -TP:

Thrm: Wasserstein Exponential Mechanism [SSR2022]<sup>1</sup>

$Y \sim f_{X, \text{WassExpMech}}$  satisfies  $\epsilon$ -TP, where, w.r.t.  $\nu_Z$ :

$$f_{X, \text{WassExpMech}}(y) \propto \exp\left(-\frac{\epsilon L_X(y)}{2\sigma(\Delta_Z)}\right)$$

where

$$\Delta_Z \triangleq \sup_{\theta_Z \in \Theta_Z} \sup_{(s_1, s_2) \in \mathbb{S}_{\text{pairs}, \text{DP}_Z}} W_\infty(\mathbb{P}(\cdot | \theta_Z, s_1), \mathbb{P}(\cdot | \theta_Z, s_2)),$$

$$\sigma(\Delta_Z) \triangleq \sup \{|L_X(y) - L_{X'}(y)| \mid x, x' \in \mathcal{X}^n, d(x, x') \leq \Delta_Z\}.$$

Extensions:

- Optimality: CLT asymptotics with additional Lipschitz regularity.
- Sampling:  $\epsilon$ -TP requires exact samples, implemented in [SRS21]<sup>2</sup>

<sup>1</sup>Seeman et al, 2022. "Formal Privacy for Partially Private Data." *Under Revision at JMLR*.

<sup>2</sup>Seeman et al, 2021. "Exact Privacy Guarantees for Sampling Algorithms Implementing the Exponential Mechanism" *NeurIPS*.

For the MVN example, we have, based on  $\theta_z = (\mu, \mu_z, \sigma^2, \Sigma_{XVZ}, \Sigma_{VZ})$

$$\Delta_z = \sup_{\theta_z \in \Theta_z} \left[ \Sigma_{XVZ}^T \Sigma_{VZ}^{-1} \begin{pmatrix} \Delta \\ z - \mu_z \end{pmatrix} \right]$$

Notes:

- $\Delta_z$  depends on *observed* value of  $z \in \mathcal{Z}$ .
- Maximum achieved by “large” (in spectral norm terms) values of  $\Sigma_{XVZ}$  and  $\Sigma_{VZ}$  (i.e., with the most dependence).

Without  $Z$ , characterizing private inference requires the marginal of  $Y$ :

$$M_{\mathcal{M},\theta}(y) \triangleq \sum_{X \in \mathcal{X}^n} \mathbb{P}_{\mathcal{M}}(y | X) \mathbb{P}_{\theta}(X)$$

Under DP, the missing data problem [RL19] has some unique properties:

- **Complete missingness:**  $X$  is not observed.
- **Perfect specification:**  $Y | X$  is *designed*, not modeled.

All disclosure avoidance methods require adjusting statistical inferences, *not just those that inject additional noise*. [SS2022]<sup>1</sup>

---

<sup>1</sup>Slavković and Seeman, 2023 “Statistical Data Privacy: A Song of Privacy and Utility.” *ARSIA*.



Two problem classes in private inference:

- **Design** an optimal privacy mechanism and estimator pair for a particular inference task under loss function  $L$ :

$$\hat{\theta}_{\text{Design}} = \arg \min_{\tilde{\theta}, \tilde{\mathcal{M}}} \sup_{\mathbb{P}_{\theta}} \mathbb{E}_{M_{\tilde{\theta}, \tilde{\mathcal{M}}}} \left[ L \left( \tilde{\theta}(Y), \theta \right) \right]$$

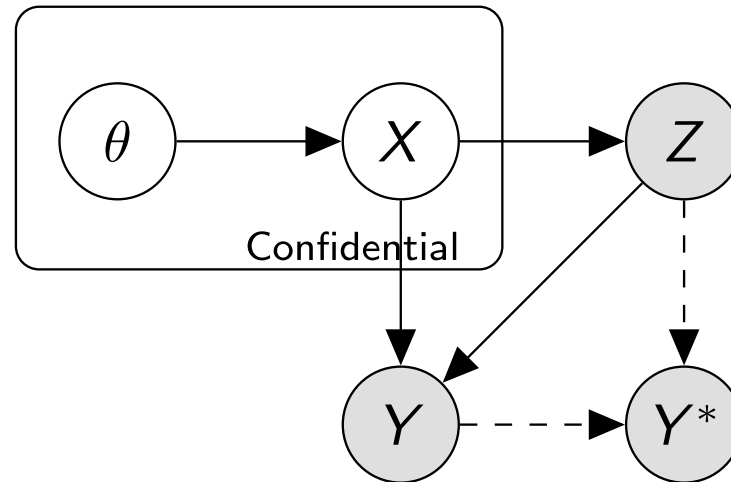
- **Adjust** an existing privacy mechanism  $\mathcal{M}$  for a particular inference task:

$$\hat{\theta}_{\text{Adjust}} = \arg \min_{\tilde{\theta}} \sup_{\mathbb{P}_{\theta}} \mathbb{E}_{M_{\tilde{\theta}, \mathcal{M}}} \left[ L \left( \tilde{\theta}(Y), \theta \right) \right]$$

Both problem classes require *adjusting uncertainty for privacy preservation*, but *optimal approaches for one may not be optimal for the other!* [SS2022]<sup>1</sup>

---

<sup>1</sup>Slavković and Seeman, 2023 “Statistical Data Privacy: A Song of Privacy and Utility.” *ARSIA*.



...however, we may not always have direct access to  $Y$  and  $Z$ !

- Private outputs may be post-processed to conform with public information, often by solving an optimization problem to minimize  $\|Y - Z\|$ .
- Ex: U.S. Census Bureau's Top-Down Algorithm [A+2021]
- **Post-processing affects design and adjustment differently in the presence of public information!**

Performing inference directly on  $Y, Z$  is preferable to using  $Y^*$ :

- 1 Finite-sample inference given  $Y, Z$  is **more powerful** than inference given  $Y^*$  in expectation, and **uniformly for exponential family models**<sup>1</sup>.
- 2 **Post-processing can remove auxiliary information** that limits inference given  $Y^*$ , even for asymptotically optimal results.
- 3 The joint distribution of  $Y, Z$  is computationally easier to work with than the marginal distribution of  $Y^*$ , and we **can derive approximate Bayesian computation (ABC) algorithms for exact sampling** from the posterior of  $\theta \mid Y, Z$  [G2019,SSR2022]<sup>2</sup>

(All theorem statements / proofs in appendix)

---

<sup>1</sup>Seeman et al, 2022. “Formal Privacy for Partially Private Data.” *Under Revision at JMLR*.

<sup>2</sup>Seeman et al, 2020. “Private Posterior Inference Consistent with Public Information: A Case Study in Small Area Estimation from Synthetic Census Data.” *PSD*.

Suppose the Pennsylvania Dept. of Health had the following release plan:

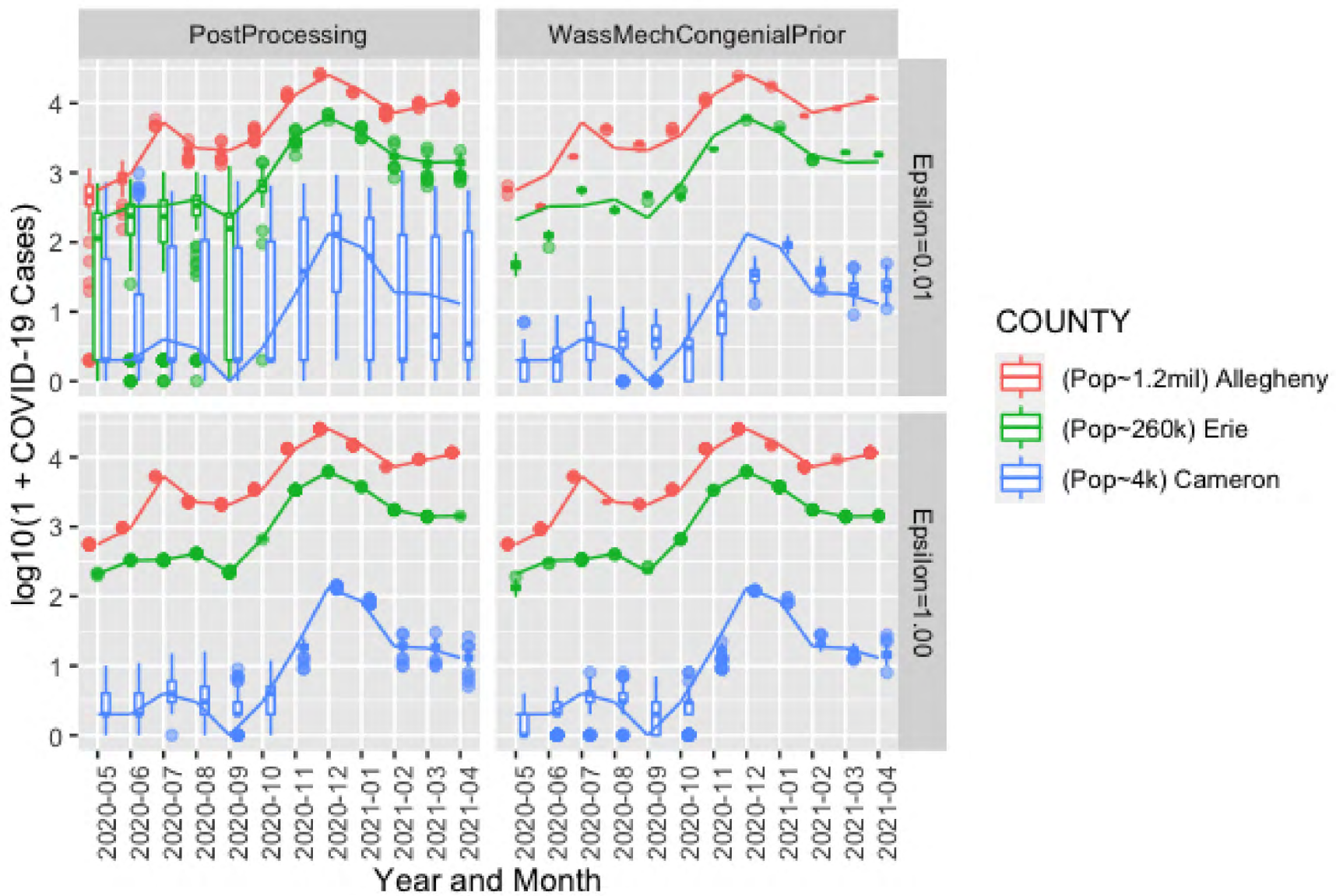
- Public statistics  $Z$ : last month's per-county COVID-19 cases and total current COVID-19 cases
- Private statistics  $Y$ : current month's per-county COVID-19 cases

We analyze two different release strategies [SSR2022]<sup>1</sup>:

- PostProcessing: calculate  $Y^*$  as the statistics closest to  $Y$  that agree with  $Z$  (all  $Y^*$  entries non-zero, sum to statewide total)
- WassMechCongenialPrior: incorporate  $Z$  through the base measure  $\nu_Z$

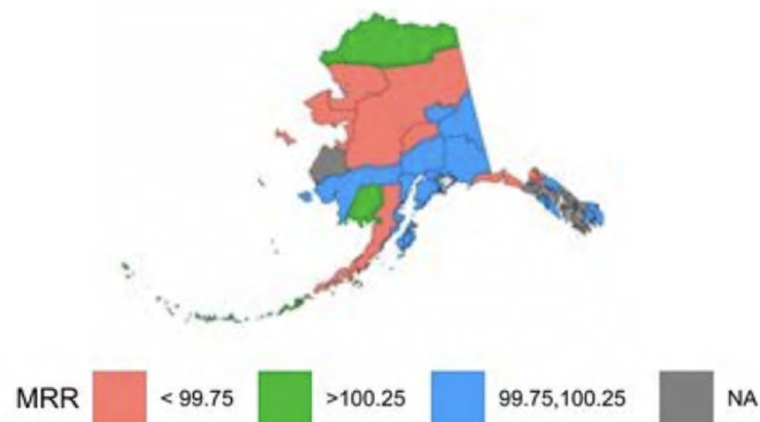
---

<sup>1</sup>Seeman et al, 2022. "Formal Privacy for Partially Private Data." *Under Revision at JMLR*.



<sup>1</sup>Seeman et al, 2022. "Formal Privacy for Partially Private Data." *Under Revision at JMLR*.

[SHV2020] demonstrated that mortality rates (using both CDC and Census data) by county have urban vs. rural and racial disparities when comparing non-private and private data released by an earlier version of the U.S. Census DP Algorithm.



**Figure:** Percentage errors in mortality rates comparing original Census private and non-private results for  $k = 14$  counties in Alaska

- Analysis outline:
  - Synthesize multiple PPD replicates according to three different methods and compare inferential accuracy
    - Private information  $X$ : small-area mortality data
    - Public information  $Z$ : nationally aggregated mortality data
- Methods to compare:
  - 1 Naive: (unadjusted) directly substitute noisy DP counts  $Y$  into test statistic calculation  $\hat{\theta}(Y)$ .
  - 2 PostProcessed: post-process  $Y, Z$  into  $Y^*$  and calculate  $\hat{\theta}(Y^*)$ .
  - 3 ConstrainedPosterior: (adjusted) estimate test statistic using empirical distribution of posterior samples drawn from  $\theta \mid Y, Z$ .

---

<sup>1</sup>Seeman et al, 2020. “Private Posterior Inference Consistent with Public Information: A Case Study in Small Area Estimation from Synthetic Census Data.” *PSD*.

For small counties in rural Alaska, sampling from  $\theta \mid Y, Z$  offers better average data utility by limiting additional errors introduced by post-processing.

County	Method	MSE	Variance	Bias
Haines	Naive	0.33	0.20	0.13
Haines	PostProcessed	0.32	0.19	0.13
Haines	<b>ConstrainedPosterior</b>	<b>0.04</b>	0.03	0.02
Nome	Naive	0.31	0.19	0.11
Nome	PostProcessed	0.31	0.20	0.11
Nome	<b>ConstrainedPosterior</b>	<b>0.03</b>	0.02	0.01
Prince of Wales	Naive	0.17	0.13	0.04
Prince of Wales	PostProcessed	0.17	0.13	0.04
Prince of Wales	<b>ConstrainedPosterior</b>	<b>0.06</b>	0.05	0.01

**Table 2.** Comparison of DP estimates of  $\hat{P}(H_1 \mid \{\vec{Y}_t\}_{t=1}^T)$  from 100 synthetic DP data sets for small counties in Alaska

<sup>1</sup>Seeman et al, 2020. “Private Posterior Inference Consistent with Public Information: A Case Study in Small Area Estimation from Synthetic Census Data.” *PSD*.



- Relaxations of DP *help account for data collection and privacy-utility trade-off decision-making*, especially for social science data like surveys.
- As more auxiliary data sets use FP, research data products like Census noisy measurements files (NMF) will provide *better inferences that require modeling noise in auxiliary data*.
- Implementing FP requires holistic risk evaluations that account for curator-induced public information for greater accountability and transparency [SS23].<sup>1</sup>

United States®  
**Census**  
Bureau



### Census Bureau Releases 2020 Census DHC Noisy Measurement File

**October 23, 2023:** The U.S. Census Bureau today released the [Noisy Measurement File](#) associated with the 2020 Census Demographic and Housing Characteristics File (DHC). The 2020 Noisy Measurement Files are considered research-based statistical products and should not be considered the official 2020 Census counts.

---

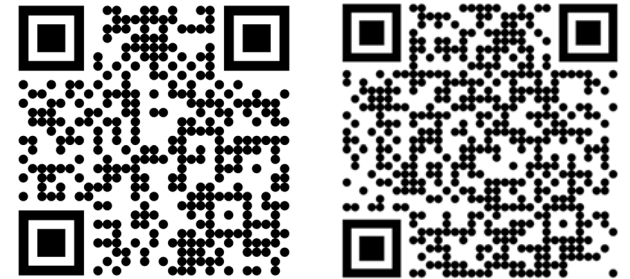
<sup>1</sup>Seeman and Susser, 2023. “Between Privacy and Utility: On Differential Privacy in Theory and Practice.” *ACM JRC*.

Contact: [jhseeman@umich.edu](mailto:jhseeman@umich.edu)

Work supported by:

- U.S. Census Bureau Dissertation Fellowship
- RA funding from NORC at UChicago and Tumult Labs
- NSF SES-1853209

Paper + Website Links:



Main reference: Seeman et al, 2022. "Formal Privacy for Partially Private Data." Under revision at JMLR. <https://arxiv.org/abs/2204.01102>

Other personal references:

- Seeman and Susser, 2023. "Between Privacy and Utility: On Differential Privacy in Theory and Practice." ACM JRC.
- Slavković and Seeman, 2023 "Statistical Data Privacy: A Song of Privacy and Utility." ARSIA.
- Seeman et al, 2021. "Exact Privacy Guarantees for Sampling Algorithms Implementing the Exponential Mechanism." NeurIPS.
- Seeman et al, 2020. "Private Posterior Inference Consistent with Public Information: A Case Study in Small Area Estimation from Synthetic Census Data." PSD.
- Seeman. "Better Privacy Theorists for Better Data Stewards." (forthcoming) JPC, 2023.

Additional references available upon request.

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

“...but I thought DP was supposed to provide protections against adversaries with arbitrary background knowledge [BGKS2012,KS2014]!? Isn't this a solved problem!?”

Short answer: not technically.

- DP ensures that any arbitrary posteriors, updated from two adjacent databases based only on the mechanism output  $Y$ , are close in statistical distance.
- DP does NOT ensure that this same relationship holds when the mechanism form changes with public information.

Let  $x_{i,v} \in \mathcal{X}^n$  be the database  $x$  where the  $i$ th record is replaced with arbitrary, data-independent value  $v$ . Define:

$$\pi_{i,v}(x | y) \triangleq \frac{\mathbb{P}(M(x_{i,v}) = y)\pi(x)}{\sum_{x^* \in \mathcal{X}^n} \mathbb{P}(M(x_{i,v}^*) = y)\pi(x^*)}.$$

Then [KS2014] show that, for any  $x, i, v, y$ , we have

$$d_{\text{TV}}(\pi(\cdot | y), \pi_{i,v}(\cdot | y)) \leq e^\epsilon - 1.$$

*Problem: what happens when  $M$  depends on  $Z$ ?*

Let  $X_1, \dots, X_n \in \{0, 1\}$  and let  $T(X) = \sum_{i=1}^n X_i$ . Then releasing  $Y$  s.t.

$$\mathbb{P}_{\epsilon, T(X)}(Y = y) = C^{-1}(\epsilon, T(X)) \exp\left(-\frac{\epsilon}{2}|T(X) - y|\right) \mathbb{1}_{\{y \in \{0, 1, \dots, n\}\}}$$

satisfies  $\epsilon$ -DP, where

$$C(\epsilon, T(X)) = \sum_{j=0}^n \exp\left(-\frac{\epsilon}{2}|T(X) - j|\right).$$

Note that  $C(\cdot, \cdot)$  is

- Monotonically decreasing as  $\epsilon$  increases.
- Monotonically decreasing as  $|T(X) - n/2|$  increases.

Suppose we want to ensure, for some  $\alpha \in \{0, 1, \dots, n\}$  s.t.  $T(X) \geq \alpha$  and  $n - T(X) \geq \alpha$  and  $\beta \in (0, 1)$ ,

$$\mathbb{P}_{\epsilon, T(X)}(|Y - T(X)| \leq \alpha) \geq 1 - \beta.$$

Equivalently,

$$C^{-1}(\epsilon, T(X)) \left( 1 + 2 \sum_{j=1}^{\alpha} \exp\left(-\frac{j\epsilon}{2}\right) \right) \geq 1 - \beta$$

Let  $\phi(\epsilon, T(X), \alpha, \beta) = 1$  if this condition is satisfied and  $\phi(\epsilon, T(X), \alpha, \beta) = 0$  otherwise. Then  $\phi(\cdot, \cdot, \cdot, \cdot)$  is

- Monotonically increasing as  $\epsilon, |T(X) - n/2|$  increases.
- Monotonically increasing in  $\alpha, \beta$  increases.



## Lemma

For any  $k \in \{\min\{T(X), n - T(X)\}, \dots, n/2\}$ , there exists  $\alpha^*, \beta^*$  such that

$$\phi(\epsilon, T(X), \alpha^*, \beta^*) = 1 \implies |T(X) - n/2| > k.$$

Suppose the PLB  $\epsilon^*$  is chosen such that  $\phi(\epsilon^*, T(X), \alpha^*, \beta^*) = 1$ . Then both the following are true:

- 1 The marginal distribution of  $Y$  satisfies the  $\epsilon^*$ -DP probability inequality under the exponential mechanism implemented with PLB  $\epsilon^*$ .
- 2 The conditional distribution of  $Y \mid \phi(\epsilon^*, T(X), \alpha^*, \beta^*) = 1$  fails to satisfy the  $\epsilon$ -DP probability inequality for any finite  $\epsilon$  value.

- Attempted solution 1: what if you set the PLB based on alternative, public data sources?
  - Ex: setting PLBs based on historical and/or synthetic data (typical approach used by Bureau, FSAs, etc.).
  - Problem: suppose the PLB  $\epsilon^*$  is chosen so that  $\phi(\epsilon^*, T(\tilde{X}), \alpha, \beta) = 1$ , where  $T(\tilde{X})$  is publicly available. Then if  $T(\tilde{X})$  and  $T(X)$  are close with high probability, the previous lemma still holds.
- Attempted solution 2: what if you chose the PLB based on the worst-case database?
  - Problem: to achieve the same data utility under an empirically chosen PLB  $\epsilon$ , a new worst-case privacy loss budget  $\epsilon^*$  takes the form
 
$$\epsilon^* = \inf \{ \epsilon' \geq \epsilon \mid \mathbb{P}_{\epsilon, T(X)}(|Y - T(X)| \geq \alpha) \leq \mathbb{P}_{\epsilon', n/2}(|Y - n/2| \leq \alpha) \}$$
  - $\wedge$  This scales linearly with the number of queries answered, when we only considered answering one sum query.

Randomized algorithms and conditional distributions are different mathematical objects without the same properties.

- DP mechanisms are collections of marginal distributions on  $Y$  indexed by realizations of the confidential data  $X$ , *which is treated as a constant*.
- Conditional distributions are *only unique up to sets of measure zero*.

## Explanation

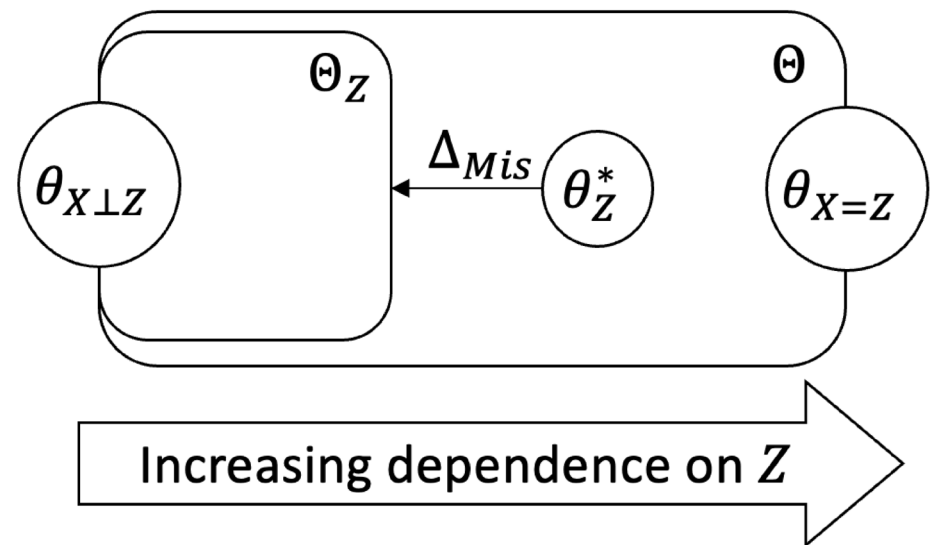
Randomized algorithms *only* characterize the presented mechanism form, **not** *how* the mechanism form was chosen and its affects on privacy guarantees.

⇒ We should analyze  $Y | Z$  to capture *how*  $Z$  changes the mechanism form.

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

How do we choose  $\Theta_Z$ , the space of admissible dependency models?

- Too little dependence: underestimating privacy losses.
- Too much dependence: failing to provide meaningful privacy guarantees (bound by “no free lunch” [K2012]).



Let  $d_\infty$  be the log-max divergence between two distributions.

### Corollary: robustness to misspecification

If the data is generated by  $\theta^* \notin \Theta_Z$ , then releasing  $Y$  satisfies  $(\epsilon + 2\Delta_{\text{mis}})$ -TP, where:

$$\Delta_{\text{mis}} = \inf_{\theta_z \in \Theta_Z} \sup_{(s_1, s_2) \in \mathbb{S}_{\text{pairs}, \text{DP}_Z}} \max \{ d_\infty(\mathbb{P}_{\theta^*, s_1}, \mathbb{P}_{\theta_z, s_1}), d_\infty(\mathbb{P}_{\theta^*, s_2}, \mathbb{P}_{\theta_z, s_2}) \}$$

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration**
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

## Many generic $\epsilon$ -DP algorithms are not exactly implementable!

Why can't we just use MCMC?

- MCMC approximation has a privacy cost
- Heuristic MCMC convergence measures tell us nothing about said cost

### Approximation $\delta$ cost (Li et al, 2016)

A sequence of mechanisms  $\mathcal{M}_m \triangleq \{\mu_{m,\mathcal{X}} \mid x \in \mathcal{X}\}$  approximating the exponential mechanism,  $\mathcal{M}_m$ , as  $m \geq \tau(\alpha)$  is  $(\epsilon, \delta_\alpha)$ -DP where  $\delta_\alpha \triangleq \alpha(1 + e^\epsilon)$  if

$$\tau(\alpha) \triangleq \sup_{\mathcal{X} \in \mathcal{X}^n} \inf \{t \geq 0 \mid \|\mu_{t,\mathcal{X}} - \mu_{\mathcal{X}}\|_{\text{TV}} \leq \alpha\}.$$



Current approach: bounding distributional distances between the MCMC approximation and the target distribution

- Total variation: for  $V$ -geometrically ergodic algorithms, there exists  $r \in (0, 1)$ ,  $C \in \mathbb{R}$  and  $V : \mathcal{Y} \mapsto \mathbb{R}^+$  s.t.:

$$\|\mu_{m,X} - \mu_X\|_{\text{TV}} \leq CV(y_0)r^m$$

Examples: Metropolis-Hastings, Hybrid MC, Hamiltonian MC

- Rényi divergence: (Ganesh and Talwar, 2020) derive asymptotic rates for Langevin dynamics step sizes and runtimes to satisfy  $(\epsilon, \alpha)$ -RDP

Problems with existing approaches:

- Asymptotic rates **can't be used to calculate finite-chain privacy loss**
- Need to bound distances for **worst-case slowest mixing chains**
- **None of the mechanisms above are  $\epsilon$ -DP**

Original approach (typical MCMC analysis):

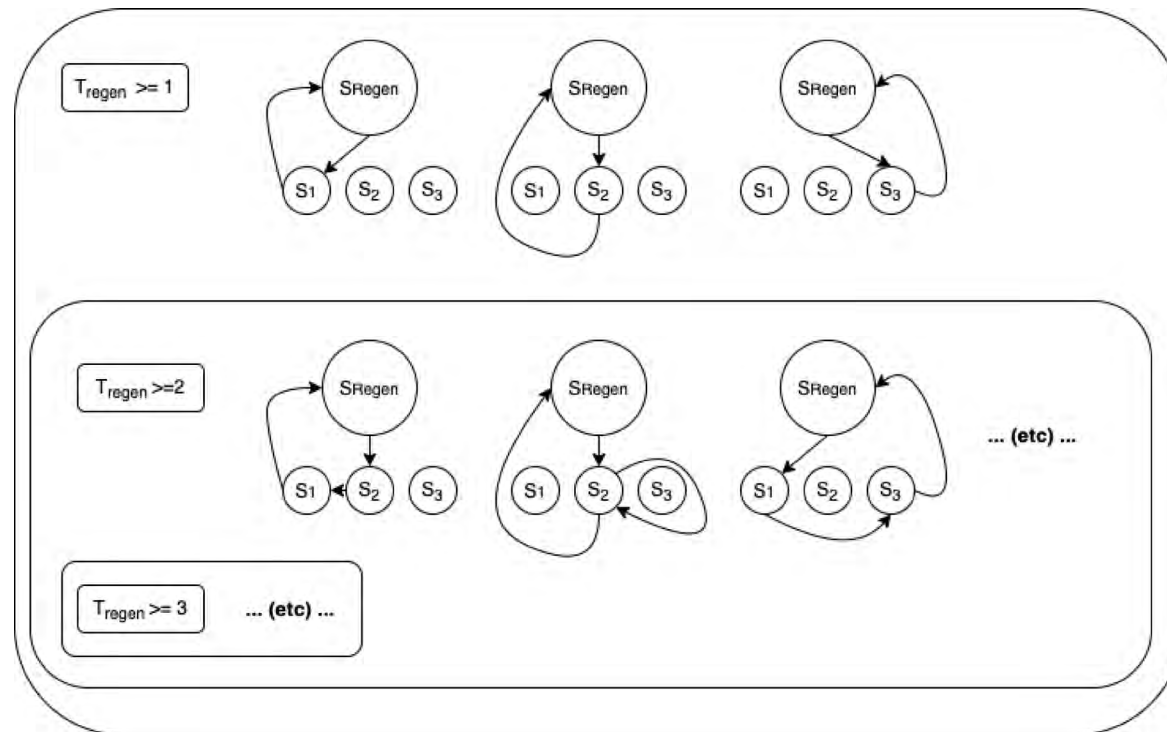
- 1 Approximate a target distribution using a guaranteed finite runtime sampling algorithm
- 2 Estimate the quality of the approximation, often by heuristic assessment of convergence

Alternative approach (perfect sampling):

- 1 Derive an algorithm that modifies an existing Markov chain with random runtime to generate an exact sample from the target distribution
- 2 Before use, determine which conditions are necessary so the expected runtime is finite and the sampling scheme is feasible

Consider a Markov chain that, w.p. 1, returns to some state  $S_{\text{regen}}$ :

- Each “tour” of the state space is IID by Markov property
- Target distribution is a mixture of the last state entered BEFORE entering  $S_{\text{regen}}$ , conditional on the times it takes to return to  $S_{\text{regen}}$



Notation:

$$\left\{ \begin{array}{l} \Pi_X : \mathcal{Y} \times \mathcal{F}_Y \mapsto [0, 1] \triangleq \text{Transition kernel} \\ Y_m : m\text{th element of original Markov chain} \\ s : \mathcal{Y} \mapsto [0, 1] \triangleq \text{Minorizing function} \\ \nu : \mathcal{F}_Y \mapsto [0, 1] \triangleq \text{Minorizing measure} \\ \xi_y \triangleq \text{Dirac delta measure at } y \in \mathcal{Y} \end{array} \right.$$

Key assumption: **minorization**

$$\Pi_X(y, A) \geq s(y)\nu(A), \quad \forall y \in \mathcal{Y}, A \in \mathcal{F}$$

Decompose  $\Pi_X$  into its minorizing kernel and a remainder kernel  $R_{\nu,s}$

$$\Pi_X(y, A) = s(y)\nu(A) + [1 - s(y)]R_{\nu,s}(y, A)$$

Next, **extend the state space** to  $\mathcal{Y} \times \{0, 1\}$  so that the random variable  $(Y, \rho)$  evolves according to:

$$\Pi_{X,\text{split}}(y, \rho; A, \rho') \triangleq [\mathbb{1}_{\{\rho=1\}}\nu(A) + \mathbb{1}_{\{\rho=0\}}R_{\nu,s}(y, A)] s(y)^{\rho'} (1 - s(y))^{1-\rho'}$$

Extended chain “regenerates” when  $\rho_m = 1$ :

$$\tau_{\nu,s} \triangleq \min\{m \in \mathbb{N} \mid \rho_{m+1} = 1\}$$

The original target distribution  $\mu_X$  is an infinite mixture conditional on regeneration times for the split chain!

$$\mu_X(A) = \sum_{m=1}^{\infty} \frac{\mathbb{P}(\tau_{\nu,s} \geq m)}{\mathbb{E}[\tau_{\nu,s}]} \mathbb{P}(Y_m \in A \mid \tau_{\nu,s} \geq m)$$

To sample from  $\mu_X$ , we need to complete two tasks in finite expected time:

- 1 Sample from  $T$  with PMF  $\mathbb{P}(T = t) = \frac{\mathbb{P}(\tau_{\nu,s} \geq t)}{\mathbb{E}[\tau_{\nu,s}]}$
- 2 Sample from  $\mathbb{P}(Y_m \in A \mid \tau_{\nu,s} \geq T)$

## Theorem (Lee et al, 2014)

Suppose there exists a singleton set  $\{a\} \in \mathcal{F}$  such that:

- 1 (Regeneration) There exists a measure  $\mu_{\text{regen}}$  on  $(\mathcal{Y}, \mathcal{F})$  such that:

$$\Pi_X(\{a\}, A) = \mu_{\text{regen}}(A) \quad \forall A \in \mathcal{F}.$$

- 2 (Uniform minorization) there exists  $\beta > 0$  such that:

$$\inf_{y \in \mathcal{Y}} \Pi_X(y, \{a\}) \geq \beta > 0$$

Then the mixture sampling method can be used to draw perfect samples from  $\mu_X$  with finite expected runtime.

Key proof techniques:

- $s(y) = \Pi_X(y, \{a\})$  and  $s(y) = \beta$  are BOTH minorizing functions w.r.t.  $\xi_a$
- Sampling from remainder using Bernoulli Factory algorithm (Huber, 2014)

Let  $(s_1, s_2) \in \mathbb{S}_{\text{pairs, DP}_z}$  and  $B \in \mathcal{F}_Y$ . Let  $\mu_{i, \theta_z} \triangleq \mathbb{P}(\cdot \mid \theta_z, s_i)$ . Let  $\gamma^*$  be the joint distribution that achieves the Wasserstein distance bound. Then

$$\frac{\mathbb{P}(Y \in B \mid s_1, \theta_z)}{\mathbb{P}(Y \in B \mid s_2, \theta_z)} = \frac{\int_{\mathcal{X}} \mathbb{P}(Y \in B \mid s_1, \theta_z, X = x) d\mu_{1, \theta_z}(x)}{\int_{\mathcal{X}} \mathbb{P}(Y \in B \mid s_2, \theta_z, X = x) d\mu_{2, \theta_z}(x)} \quad (1)$$

$$= \frac{\int_{\mathcal{Y}} \int_{\mathcal{X}} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_x(y)}{2\sigma(\Delta_z)}\right) d\mu_{1, \theta_z}(x) d\nu_z(y)}{\int_{\mathcal{Y}} \int_{\mathcal{X}} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_x(y)}{2\sigma(\Delta_z)}\right) d\mu_{2, \theta_z}(x) d\nu_z(y)} \quad (2)$$

$$= \frac{\int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_x(y)}{2\sigma(\Delta_z)}\right) d\gamma^*(x, x') d\nu_z(y)}{\int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{\mathcal{X}} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_{x'}(y)}{2\sigma(\Delta_z)}\right) d\gamma^*(x, x') d\nu_z(y)}. \quad (3)$$



Let  $B_{\Delta_z}(x) \triangleq \{x' \in \mathcal{X} \mid d(x, x') \leq \Delta_z\}$ . Then by construction and definition of the Wasserstein distance,

$$\frac{\mathbb{P}(Y \in B \mid s_1, \theta_z)}{\mathbb{P}(Y \in B \mid s_2, \theta_z)} \quad (4)$$

$$= \frac{\int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{B_{\Delta_z}(x)} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_x(y)}{2\sigma(\Delta_z)}\right) d\gamma^*(x, x') d\nu_z(y)}{\int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{B_{\Delta_z}(x')} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_{x'}(y)}{2\sigma(\Delta_z)}\right) d\gamma^*(x, x') d\nu_z(y)} \quad (5)$$

$$\leq \frac{\int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{B_{\Delta_z}(x')} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon(L_{x'}(y) + \sigma(\Delta_z))}{2\sigma(\Delta_z)}\right) d\gamma^*(x, x') d\nu_z(y)}{\int_{\mathcal{Y}} \int_{\mathcal{X}} \int_{B_{\Delta_z}(x')} \mathbb{1}_{\{y \in B\}} \exp\left(-\frac{\epsilon L_{x'}(y)}{2\sigma(\Delta_z)}\right) d\gamma^*(x, x') d\nu_z(y)} \quad (6)$$

$$\leq \exp\left(\frac{\epsilon}{2}\right) \exp\left[\frac{\epsilon}{2} \left(\frac{\sup\{|L_x(y) - L_{x'}(y)| \mid x, x' \in \mathcal{X}^n, d(x, x') \leq \Delta_z\}}{\sigma(\Delta_z)}\right)\right] \quad (7)$$

$$\leq \exp(\epsilon). \quad (8)$$

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG**
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

Regularity conditions:

- 1  $n^{-1}\nabla^2 L_{X_n}(y)$  exists almost everywhere  $\nu_{Z_n}$ , and the smallest eigenvalue of  $L_{X_n}(y)$  is greater than  $\alpha > 0$  for all  $n \in \mathbb{N}$ ,  $y \in \mathcal{Y}$ ,  $\nu_{Z_n}$  almost everywhere.
- 2 There exists a unique  $y^* \in \mathcal{Y}$  such that as  $n \rightarrow \infty$ ,

$$Y_n^* \triangleq \arg \min_{\tilde{y} \in \mathcal{Y}} L_{X_n}(\tilde{y}) \rightarrow_P y^*. \quad (9)$$

- 3 For all  $n \in \mathbb{N}$ ,  $\Delta_{Z_n} \geq \Delta^* > 0$  w.p. 1
- 4  $\sigma_y(\Delta^*)$  is continuous in  $y$ , and  $\sigma_y(\Delta^*) \geq \sigma^*(\Delta^*)$  holds  $\nu_{Z_n}$  almost everywhere.

Regularity conditions (continued):

5 For all  $n \in \mathbb{N}$ ,

$$\int_{\mathcal{X}} \exp\left(-\frac{\alpha \|y_n - y_n^*\|_K}{2\sigma_{y_n^*}(\Delta^*)}\right) d\nu_{Z_n} < \infty. \quad (10)$$

6  $Z_n \rightarrow_D Z^*$  and  $\nu_{Z_n} \rightarrow_D \nu_{Z^*} \ll \lambda$ .

7 Let  $B_\delta(\cdot)$  be a  $K$ -norm ball of radius  $\delta$ . We assume that for  $a \in \mathcal{Y}$  and some  $\delta > 0$ ,

$$\int_{B_\delta(y_n^*)} d\nu_{Z_n}(a) = \Omega_P(1). \quad (11)$$

In words, the base measure should support the true solution  $y^*$  with probability bounded away from zero by some constant.

Let  $f_n(y)$  be the density of the Wasserstein  $K$ -norm gradient mechanism,

$$f_n(y_n) \propto \exp\left(-\frac{\epsilon \|\nabla L_{X_n}(y_n)\|_K}{2\sigma_y(\Delta_{Z_n})}\right) d\nu_{Z_n}(y_n). \quad (12)$$

Let  $a_n \triangleq n(y_n - y_n^*)$  with density  $h_n(\cdot)$  so that

$$h_n(a_n) \propto \exp\left(-\frac{\epsilon \|\nabla L_{X_n}(y_n^* + a_n/n)\|_K}{2\sigma_y(\Delta_{Z_n})}\right) d\nu_{Z_n}(y_n^* + a_n/n). \quad (13)$$

Note that in the transformation above, the Jacobian constant  $n^{-1}$  gets absorbed in the proportionality.

Next, we Taylor expand the loss function,

$$\nabla L_{X_n}(y_n) = \nabla L_{X_n}(y_n^* + a_n/n) = \nabla^2 L_{X_n}(y_n^*) \frac{a_n}{n} + o_p(1) \quad (14)$$

Following the proof of [RA2019], using Assumptions (1) - (4), there exists a constant  $C > 0$  such that

$$-\frac{1}{\sigma_{y_n^* + a_n/n}(\Delta_{Z_n})} \|\nabla L_n(y_n^* + a_n/n)\|_K \leq -\frac{C\alpha}{\sigma^*(\Delta^*)} \|a_n\|_2. \quad (15)$$

Then by the dominated convergence theorem, Assumptions (5) - (6), and the continuity of  $\sigma$ , we conclude that

$$\lim_{n \rightarrow \infty} h_n(a_n) = h(a) \propto \exp\left(-\frac{\epsilon}{2\sigma_{y^*}(\Delta^*)} \|\Sigma^{-1}a\|_K\right), \quad (16)$$

with respect to the base measure  $\nu_{Z^*}(a)$  (note that Assumption (6) ensures the asymptotic base measure supports  $y^*$ ).

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime**
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

## Theorem: modified Metropolis-Hastings perfect sampling for privacy

Let  $\Pi_X$  be the transition kernel for a Metropolis-Hastings Markov Chain with symmetric proposals  $q$ . We can construct a Markov chain on the extended space with proposals:

$$\tilde{q}(y, y') = \frac{1}{2} [q_X(y, y') + \mathbb{1}_{\{y'=a\}}],$$

And an algorithm to sample from density  $f_X$  that satisfies  $\epsilon$ -DP with expected number of total proposed samples  $N_{\text{prop}}$ :

$$\mathbb{E}[N_{\text{prop}}] \leq \frac{48}{k^2(1-k)^2 \inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y)},$$

where:

$$p_{\text{Accept}}(y) \triangleq \int_{\mathcal{Y}} q_X(y, y') \min \left\{ 1, \frac{f_X(y')}{f_X(y)} \right\} d\nu(y').$$



Key ingredient: choose  $s(x) \triangleq s \in (0, \beta)$  and  $\nu = \xi_a$  (point mass at  $a$ )

$$\mu_X(A) = \sum_{m=1}^{\infty} s(1-s)^{m-1} \left[ \xi_a R_{\xi_a, s}^{m-1} \right]$$

To sample from  $\mu_X$ , we need:

- 1 Sample  $T \sim \text{Geometric}(s)$
- 2 Sample from  $Y \mid T = t \sim \xi_a R_{\xi_a, s}^{t-1}$

Residual kernel has a special mixture form:

$$R_{\xi_a, s}(y, A) = \frac{1 - \Pi_X(y, \{a\})}{1 - s} R_{\xi_a, \Pi_X(\cdot, \{a\})}(y, A) + \frac{\Pi_X(y, \{a\}) - s}{1 - s} \xi_a(A)$$

Each mixture is easy to sample, but need to choose a component using:

$$\text{Bernoulli} \left( \frac{1 - \Pi_X(y, \{a\})}{1 - s} \right)$$

This is called a **Bernoulli factory problem** (Huber, 2014)

- “Using a black box simulating  $\text{Bernoulli}(\theta)$ , simulate  $\text{Bernoulli}(f(\theta))$  for known  $f$  and unknown  $\theta$ ”
- (Huber, 2014) derives efficient algorithms for solving this problem

$$\left\{ \begin{array}{l} N_{\text{Outer}} \triangleq \text{Number of outer loops in main algorithm} \\ N_{\text{Inner}} \triangleq \text{Number of rejection proposals to sample from } Y_m^* \sim R_{\xi_a, \tilde{\pi}_{\cdot, \{a\}}}(Y_{m-1}, \cdot) \\ N_{\text{Bern}} \triangleq \text{Number of Bernoulli factory } p\text{-flips to select } f(p) \text{ flip} \\ N_{\text{Nonatomic}} \triangleq \text{Number of samples from } \tilde{\mu}_X \text{ required to sample from } \mu_X \end{array} \right.$$

$$\left\{ \begin{array}{l} \mathbb{E}[N_{\text{Outer}}] = \frac{2}{k} \\ \mathbb{E}[N_{\text{Inner}}] \leq ((1 - k) \inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y))^{-1} \\ \mathbb{E}[N_{\text{Bern}}] \leq \frac{24}{k} \\ \mathbb{E}[N_{\text{Nonatomic}}] = (1 - k)^{-1} \end{array} \right.$$

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance**
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

## Theorem

Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_\theta$  for an exponential family with sufficient statistic  $T(\cdot)$ . Let  $Y$  be an instance of WEM with density

$$g_x(y) \propto \exp\left(-\frac{\epsilon}{2\sigma(\Delta_z)} L(\|y - T(x)\|)\right), \quad (17)$$

for some loss function  $L$  that depends only on a norm  $\|y - T(x)\|$ . Let  $Z = h(X)$  and let  $Y^* = \text{Proj}(Y, \mathcal{Z})$ , so that  $Y, Y^* \in \mathcal{Y}$ . Then  $Y$  has a monotone likelihood ratio in  $\theta$ . Furthermore, define the test:

$$H_0 : \theta \leq \theta_0, \quad H_1 : \theta > \theta_0. \quad (18)$$

For any unbiased test  $\phi : \mathcal{Y}^* \mapsto [0, 1]$  for  $\theta$  based on  $Y^*$ , there exists a uniformly more powerful test  $\phi' : \mathcal{Y} \times \mathcal{Z} \mapsto [0, 1]$ . If  $\mathbb{P}_\theta(Y \neq Y^*) > 0$  for all  $\theta \in \Theta$ , then this improvement is strict.

## Main points:

- 1 Using [W1985], the marginal of  $Y$  has MLR in  $\mathcal{T}$ .
- 2 If  $Y^*$  is not a bijection, multiple  $Y$ s get mapped to the same  $Y^*$  which allows for a tighter bound on Type I error.
- 3  $\uparrow$  This can be used to construct the uniformly more powerful test.

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling**
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

**Result:**  $Y \sim \tilde{\mu}_X$

**Input:** Transition kernel  $\Pi_X$ , singleton atom  $a \in \mathcal{Y}$ , minorization  $\beta > 0$

Sample  $M \sim \text{Geometric}(\beta)$ , set  $Y_1 = a$  ;

**for**  $m = 2$  **to**  $M$  **do**

Using the Bernoulli factory algorithm [H2016], sample:

$$Z_m \sim \text{Bernoulli} \left( \frac{1 - \tilde{\Pi}_X(Y_{m-1}, \{a\})}{1 - \beta} \right). \quad (19)$$

**if**  $Z_m = 1$  **then**

Sample  $Y_m \sim R_{\xi_a, \tilde{\Pi}_X(\cdot, \{a\})}(Y_{m-1}, \cdot)$  using rejection:

- 1 Propose  $Y_m^* \sim \tilde{\Pi}_X(Y_{m-1}, \cdot)$ .
- 2 Accept if  $Y_m^* \neq a$ , else go back to 1.

**else**

|  $Y_m = a$

**end**

Release  $Y_M$ .

**end**

**Algorithm 1:** Exact implementation for sampling from a singleton atom mixture



Problem: singleton atoms don't usually exist

Solution: make our own! (Brockwell and Kadane, 2006)

- 1 Extend the state space with an artificial atom at  $a \in \mathcal{Y}$  to yield a new target density:

$$\tilde{\mu}_X(A) \triangleq (1 - k)\mu_X(A) + k\mathbb{1}_{\{a \in A\}}$$

w.r.t. a new base measure:

$$\tilde{\nu}(A) = (1 - k)\nu(A) + k\xi_a$$

- 2 Modify the transition kernel:

$$\tilde{\Pi}_X(x, A) = w\Pi_X(x, A) + (1 - w)\Pi'_X(x, A)$$

Where  $\Pi'_X$  transitions between the two mixture components

- 3 Apply the previous algorithm to sample from  $\tilde{\mu}_X$
- 4 Use samples from  $\tilde{\mu}_X$  to sample from  $\mu_X$

Implementation choices specific to DP:

- Choose  $a \in \mathcal{Y}$  from the set of confidential results (i.e. what we would release without privacy preservation)

$$a \in \arg \inf_{y \in \mathcal{Y}} L_X(y)$$

- Assumptions about the state space (such as compact  $\mathcal{X}^n$ ) help to satisfy our privacy AND our sampling assumptions:

$$\begin{cases} \exists \Delta_L < \infty \text{ s.t. } |L_X(y) - L_{X'}(y)| \leq \Delta_L \\ \exists \beta > 0 \text{ s.t. } \inf_{y \in \mathcal{Y}} \tilde{\Pi}_X(y, \{a\}) \geq \beta \end{cases}$$

- Many different possible choices for chain modification (ex: Brockwell and Kadane, 2005)
  - Need to balance time spent in the artificial atom and original chain
  - Properties of  $\Pi_X$  don't necessarily translate to  $\tilde{\Pi}_X$

**Result:**  $Y \sim f_x$

**Input:** Sample space  $\mathcal{Y}$  and loss function  $L_x$ .

**while** TRUE **do**

    Sample  $\tilde{Y} \sim \tilde{f}_x$  using Algorithm from [L+2014].

**if**  $\tilde{Y} \neq a$  **then**

        | Release  $\tilde{Y}$ .

**end**

**end**

**Algorithm 2:** ConfAtomPerfect:  $\epsilon$ -DP exact sample from exponential mechanism

To additionally privatize runtime, we need the number of inner-loop proposals to be 0-DP. We assume an adversary knows a modified  $\tilde{N}_{\text{Inner}}$  where

$$\tilde{N}_{\text{Inner}} \triangleq \rho N_{\text{Inner}} + (1 - \rho)(N_{\text{Inner}} + N_{\text{Wait}}), \quad (20)$$

where,

$$\rho \sim \text{Bernoulli} \left( \frac{\inf_{x \in \mathcal{X}^n} \inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y)}{\inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y)} \right), \quad (21)$$

$$N_{\text{Wait}} \sim \text{Geometric} \left( \inf_{x \in \mathcal{X}^n} \inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y) \right). \quad (22)$$

Sample  $M \sim \text{Geometric}(\beta)$ , set  $Y_1 = a$  ;

**for**  $m = 2$  **to**  $M$  **do**

  : [Same as previous algorithm].

**if**  $Z_m = 1$  **then**

    : [Same as previous algorithm].

    Sample

$$\rho \sim \text{Bernoulli} \left( \frac{\inf_{x \in \mathcal{X}} \inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y)}{\inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y)} \right), \quad (23)$$

**if**  $\rho = 1$  **then**

      Set

$$\tilde{n}_{\text{inner}} = n_{\text{inner}}. \quad (24)$$

**else**

      Sample

$$n_{\text{wait}} \sim \text{Geometric} \left( \inf_{x \in \mathcal{X}} \inf_{y \in \mathcal{Y}} p_{\text{Accept}}(y) \right). \quad (25)$$

      Update:

$$\tilde{n}_{\text{inner}} \triangleq n_{\text{inner}} + n_{\text{wait}}. \quad (26)$$

**end**

**else**

$Y_m = a$

**end**

  Release  $Y_M$ .

**end**

**Algorithm 3:** Modification to [L+2014] to privatize runtime.

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference**
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

**Result:** One sample from  $\theta \mid Y, Z$

Sample  $\theta^* \sim \pi(\theta \mid Z)$  ;

Sample  $X^* \sim \pi(\cdot \mid \theta^*, Z)$  ;

Sample  $U \sim \text{Unif}(0, 1)$  ;

**if**

$$\frac{\pi(Y \mid X^*, Z)}{\sup_{y^* \in \mathcal{Y}} \pi(y^* \mid X^*, Z)} \leq U, \quad (27)$$

**then**

| Return  $\theta^*$ .

**else**

| Go to beginning.

**end**

**Algorithm 4:** Posterior rejection sampling conditional on public information

**Result:** Estimate of  $\hat{\mathbb{E}}[a(\theta) \mid Y, Z]$  using proposals from density  $g$

Sample  $\theta^{(1)}, \dots, \theta^{(m)} \sim \pi(\theta \mid Z = z)$  ;

Sample  $X^{(j)} \sim \pi(\cdot \mid \Theta = \theta^{(j)}, Z = z)$  ;

Calculate

$$w^{(j)} = \frac{\pi(Y \mid X^{(j)}, Z)\pi(\theta^{(j)} \mid Z)}{g(\theta^{(j)})}. \quad (28)$$

Return

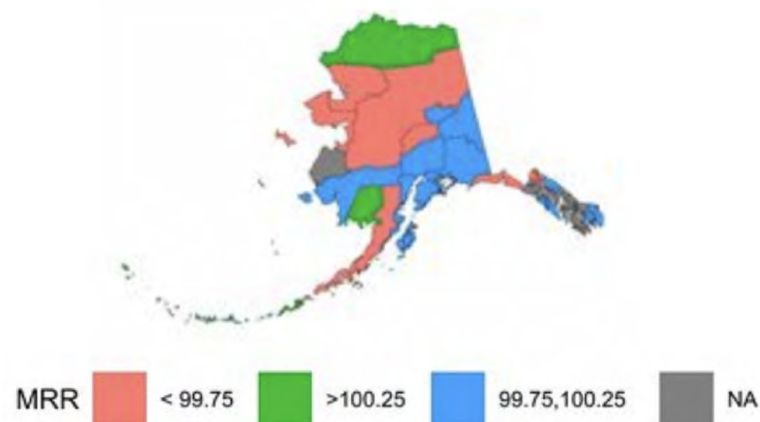
$$\hat{\mathbb{E}}[a(\theta) \mid Y, Z] \triangleq \frac{\sum_{j=1}^m w^{(j)} a(\theta^{(j)})}{\sum_{j=1}^m w^{(j)}}. \quad (29)$$

**Algorithm 5:** Posterior importance sampling conditional on public information



- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaksa mortality**
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

(Santos-Lozada et al, 2020) demonstrated that mortality rates (using both CDC and Census data) by county have urban vs. rural and racial disparities when comparing non-private and private data released by an earlier version of the U.S. Census DP Algorithm.



**Figure:** Percentage errors in mortality rates comparing original Census private and non-private results for  $k = 14$  counties in Alaska

Public information for contingency tables is often restricted to linear constraints of the form:  $C_j(\vec{Y}) = \mathbb{1}\{Q_j \vec{Y} \text{ op}_j \vec{c}_j\}$  for  $j \in [J]$

$$\begin{cases} Q_j & \text{Constraint query matrix} \\ \text{op}_j & \text{Elementwise comparison operator} \\ \vec{c}_j & \text{Constraint value} \end{cases}$$

For the non-private data,  $P(C_j(\vec{X}) = 1 \quad \forall j \in [J]) = 1$ .

- Ex1: structural inequalities:  $X_j \leq X_k$
- Ex2: exact marginals:  $\vec{a}^T \vec{X} = n_a$ .

Post-processing takes  $\vec{Y}$  as-is and solves the optimization problem:

$$\vec{Z} \triangleq h(\vec{Y}) = \arg \min_{\vec{Z}^*} \|\vec{Z}^* - \vec{Y}\|_2^2 \quad \text{s.t.} \quad Q_j \vec{Z}^* \text{ op}_j \vec{c}_j, \quad \forall j \in [J]$$

- Instead of post-processing to ensure constraints are verified, **incorporate constraint information into data generating process**
- Exploit separability of different sample generation stages:
  - Ex1: structural inequalities:  $X_j \leq X_k \implies P(\theta_j \leq \theta_k) = 1$
  - Ex2: exact marginals:  $\vec{a}^T \vec{X} = n_a \implies P(\vec{a}^T \vec{\varepsilon} = \vec{a}^T \vec{Y}_{\text{obs}} - n_a) = 1.$
- Define **model-dependent** constraints  $\mathcal{M}$  and **data-dependent** constraints  $\mathcal{D}$  as applied to  $\vec{\theta}$  and  $\vec{\varepsilon}$ :

$$C_{\mathcal{M}}(\vec{\theta}) \triangleq \prod_{m \in \mathcal{M}} m(\vec{\theta}), \quad C_{\mathcal{D}}(\vec{\varepsilon}) \triangleq \prod_{d \in \mathcal{D}} d(\vec{\varepsilon})$$

**Data:** Observed DP data  $\vec{Y}_{\text{obs}}$ , model-dependent constraints  $\mathcal{M}$ , conditional prior  $\pi_{\vec{\theta} | C_{\mathcal{M}}(\vec{\theta})=1}$ , data-dependent constraints  $\mathcal{D}$ , likelihood  $f_{\vec{X} | \vec{\theta}}$ , conditional error density  $g_{\vec{\varepsilon} | C_{\mathcal{D}}(\vec{\varepsilon})=1}$ ,

**Result:**  $N$  samples from  $\vec{\theta} | \vec{Y}$

**while**  $i \leq N$  **do**

Sample  $\vec{\theta}^{(i)} \sim \pi(\cdot | C_{\mathcal{M}}(\vec{\theta}) = 1)$ ,  $\vec{X}^{(i)} | \vec{\theta}^{(i)} \sim f(\cdot | \vec{\theta}^{(i)})$ ,  $U \sim \text{Unif}(0, 1)$  ;

**if**

$$U \leq \frac{g(\vec{\varepsilon} = \vec{Y}_{\text{obs}} - \vec{X}^{(i)} | C_{\mathcal{D}}(\vec{\varepsilon}) = 1)}{\sup_{\vec{Y} \in \Omega_Y} g(\vec{\varepsilon} = \vec{Y}_{\text{obs}} - \vec{X}^{(i)} | C_{\mathcal{D}}(\vec{\varepsilon}) = 1)}$$

**then**

| Accept sample  $\vec{\theta}^{(i)}$ ,  $i \mapsto i + 1$  ;

**else**

| Reject sample  $\vec{\theta}^{(i)}$  ;

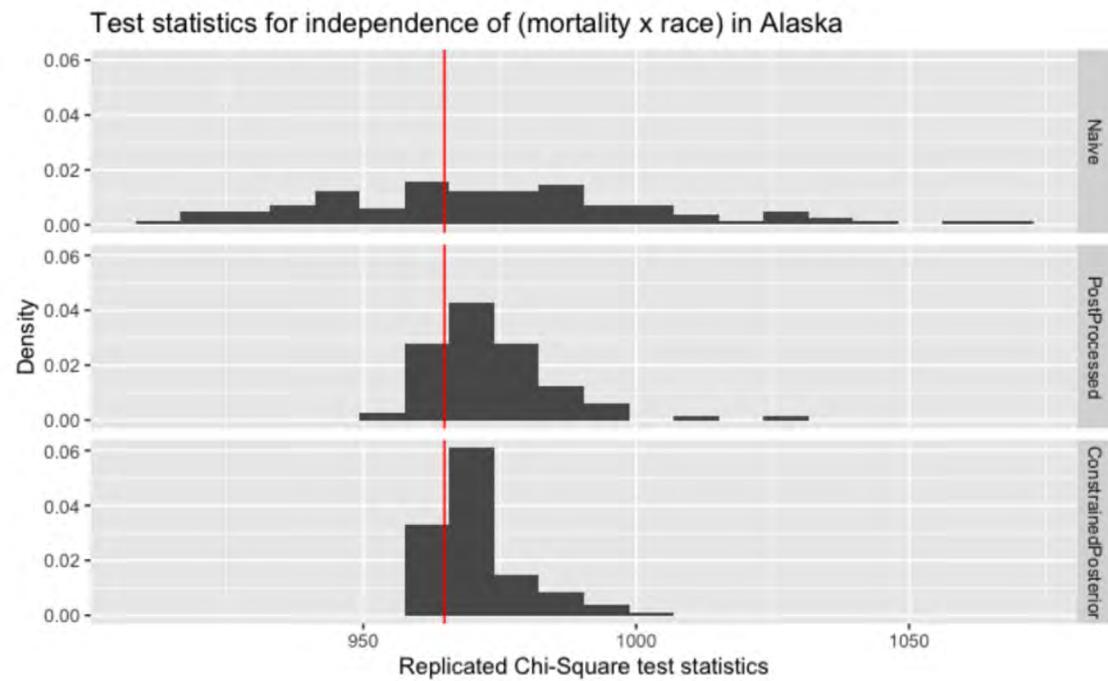
**end**

**end**

**Algorithm 6:** Constrained posterior method: DP posterior sampling given additive perturbation and public knowledge constraints

- Analysis outline:
  - Generate multiple copies of private synthetic data, and post-process each copy according to given public information
  - Compare accuracy of inferential test statistics across each synthetic data replicate
- Methods to compare:
  - 1 Naive: directly substitute noisy DP counts  $\vec{Y}$  into test statistic calculation
  - 2 PostProcessed: directly substitute post-processed DP counts  $\vec{Z}$  into test statistic calculation
  - 3 ConstrainedPosterior: estimate test statistic using empirical distribution of constrained posterior samples

- $H_0$ : mortality rate by race at national level equals mortality rate by race at state level
- Test statistic:  $\chi_{\text{obs}}^2$
- Public information: state deaths, state total population
- Constraints:
  - County level deaths between 0 and state total
  - County level population between 0 and state total
  - County level deaths smaller than county level population
  - State level marginal deaths agree with public data
  - State level marginal population agrees with public data



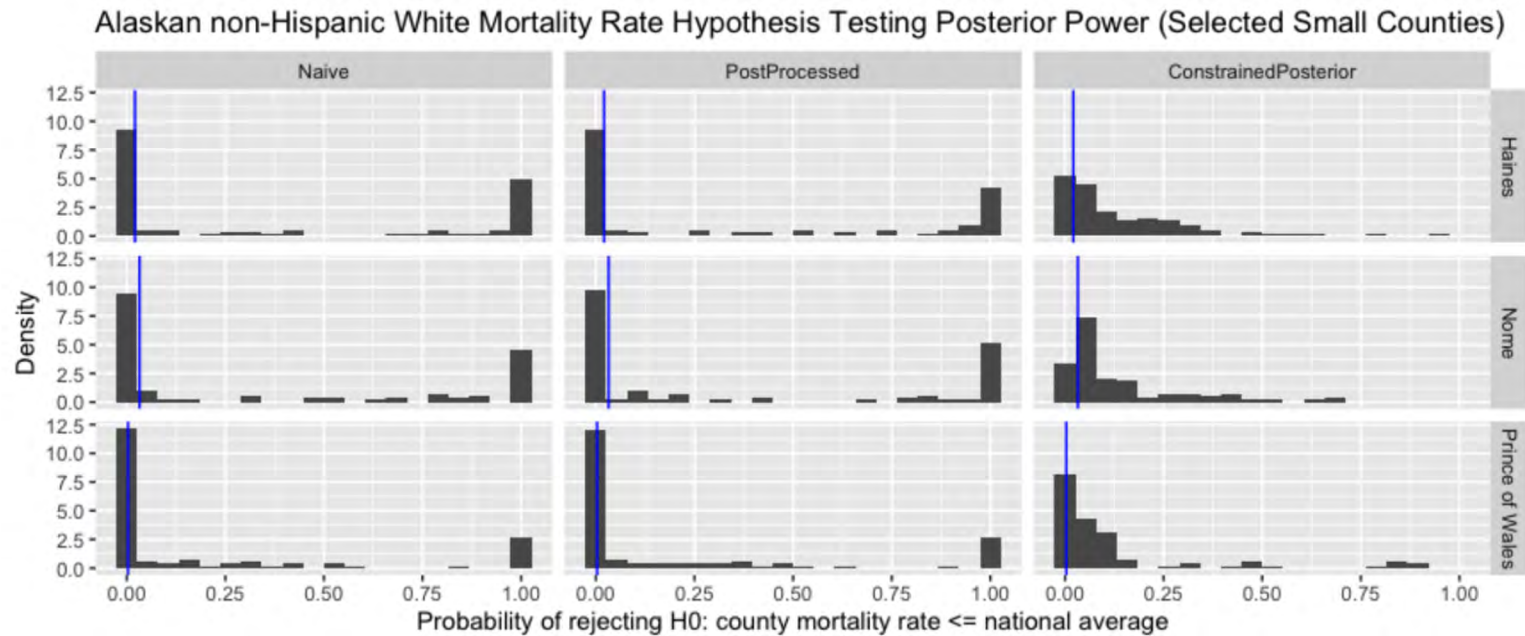
**Fig. 2.** Distributions of DP estimates of  $\hat{\chi}_{\text{obs}}^2$  test statistic from 100 synthetic DP data sets (true  $\chi_{\text{obs}}^2$  test statistic in red)



ind	SampleMSE	SampleVar	SampleBias2
Naive	1053	986	77
PostProcessed	201	130	72
ConstrainedPosterior	102	62	41

**Table 1.** Comparison of DP estimates of  $\hat{\chi}_{\text{obs}}^2$  test statistic from 100 synthetic DP data sets (true value  $\chi_{\text{obs}}^2 \approx 964$  on  $(2 - 1) \times (3 - 1)$  degrees of freedom (two mortality statuses by 3 race groups))

- $H_0$ : White mortality rate at national level equals white mortality rate by race at county level
- Test statistic:  $P(H_1 | \vec{X}_{\text{obs}})$
- Public information: county level white population
- Constraints:
  - County level white deaths between 0 and county level white population
  - County level marginal populations agree with public data



**Fig. 3.** Posterior distributions of  $\hat{P}(H_1 \mid \{\vec{Y}_t\}_{t=1}^T)$  using the three different methods outlined above. Blue lines represent the true probabilities for each county.

County	Method	MSE	Variance	Bias
Haines	Naive	0.33	0.20	0.13
Haines	PostProcessed	0.32	0.19	0.13
Haines	<b>ConstrainedPosterior</b>	<b>0.04</b>	0.03	0.02
Nome	Naive	0.31	0.19	0.11
Nome	PostProcessed	0.31	0.20	0.11
Nome	<b>ConstrainedPosterior</b>	<b>0.03</b>	0.02	0.01
Prince of Wales	Naive	0.17	0.13	0.04
Prince of Wales	PostProcessed	0.17	0.13	0.04
Prince of Wales	<b>ConstrainedPosterior</b>	<b>0.06</b>	0.05	0.01

**Table 2.** Comparison of DP estimates of  $\hat{P}(H_1 | \{\vec{Y}_t\}_{t=1}^T)$  from 100 synthetic DP data sets for small counties in Alaska

- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data**
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

## Data (available through PA Dept. of Health and IPUMS)

- Private data:
  - County-level COVID-19 case rates at month  $t$
  - County-level COVID-19 death rates at month  $t$
- Public data:
  - State-level COVID-19 case rate at month  $t$
  - County-level COVID-19 case rates at time  $t - 1$

## Synthesis methods:

- 1 Geometric noise with  $L_2$ - $L_1$  post-processing for congeniality
- 2 WassExpMech with three different base measures:
  - Improper uniform over integers
  - Congenial with public state-level case rate
  - $\uparrow +$  Dirichlet prior on  $t - 1$  county rates

With  $[n] \triangleq \{1, 2, \dots, n\}$ , define:

$$\begin{cases} j \in \{1, \dots, J\} & \triangleq \text{Pennsylvania counties, } J = 67 \\ t \in \{1, \dots, T\} & \triangleq \text{Year-month periods, } T = 24 \\ X_{j,t}^{(c)} & \triangleq \text{Number of COVID-19 cases in county } j \text{ at time } t \\ X_{j,t}^{(d)} & \triangleq \text{Number of COVID-19 deaths in county } j \text{ at time } t \end{cases} \quad (30)$$

With public information:

$$Z_t = \begin{cases} X_{j,t-1}^{(c)} = x_{j,t-1}^{(c)} \\ X_{j,t-1}^{(d)} = x_{j,t-1}^{(d)} \\ \sum_{j=1}^J X_{j,t} = s_{j,t}^{(c)} \\ \mathbb{P}(X_{j,t}^{(c)} \geq X_{j,t}^{(d)}) = 1. \end{cases} \quad (31)$$

First, we synthesize

$$\begin{cases} Y_{j,t}^{(c)} = X_{j,t}^{(c)} + \varepsilon_{j,t}^{(c)} \\ Y_{j,t}^{(d)} = X_{j,t}^{(d)} + \varepsilon_{j,t}^{(d)}, \end{cases} \quad (32)$$

where

$$\varepsilon_{j,t}^{(c)}, \varepsilon_{j,t}^{(d)} \stackrel{\text{iid}}{\sim} \text{DiscreteLaplace} \left( \frac{\epsilon}{\Delta} \right). \quad (33)$$

Then, we will perform deterministic two-stage post-processing; first, we will find the solution to the  $L_2$  optimization problem,

$$\begin{pmatrix} \tilde{Y}_t^{(c)} \\ \tilde{Y}_t^{(d)} \end{pmatrix} = \arg \min_{Y \in \mathbb{R}^{2J}} \left\| \begin{pmatrix} Y^{(c)} \\ Y^{(d)} \end{pmatrix} - \begin{pmatrix} Y_t^{(c)} \\ Y_t^{(d)} \end{pmatrix} \right\|_{L_2}^2 \quad \text{s.t.} \quad (34)$$

$$\begin{cases} \sum_{j=1}^J y_j^{(c)} = s_{j,t} \\ y_j^{(c)} \geq y_j^{(d)} \geq 0 \quad \forall j \in [J]. \end{cases} \quad (35)$$



Next, we will find the solution to the integer  $L_1$  optimization problem,

$$\begin{pmatrix} Y_t^{(c)*} \\ Y_t^{(d)*} \end{pmatrix} = \begin{pmatrix} \lfloor \tilde{Y}_t^{(c)} \rfloor \\ \lfloor \tilde{Y}_t^{(d)} \rfloor \end{pmatrix} + \arg \min_{Y \in \{0,1\}^{2J}} \left\| \begin{pmatrix} Y^{(c)} \\ Y^{(d)} \end{pmatrix} - \begin{pmatrix} \tilde{Y}_t^{(c)} \\ \tilde{Y}_t^{(d)} \end{pmatrix} \right\|_{L_1} \text{ s.t.} \quad (36)$$

$$\begin{cases} \sum_{j=1}^J y_j^{(c)} = s_{j,t} - \lfloor \tilde{Y}_t^{(c)} \rfloor \\ y_j^{(c)} + \lfloor \tilde{Y}_t^{(c)} \rfloor \geq y_j^{(d)} \lfloor \tilde{Y}_t^{(d)} \rfloor \geq 0 \quad \forall j \in [J]. \end{cases} \quad (37)$$

- 1 Naive base measure:

$$\nu_Z(Y) \propto \mathbb{1}_{\{Y \in \mathbb{Z}^J\}}. \quad (38)$$

- 2 Deterministic congenial base measure:

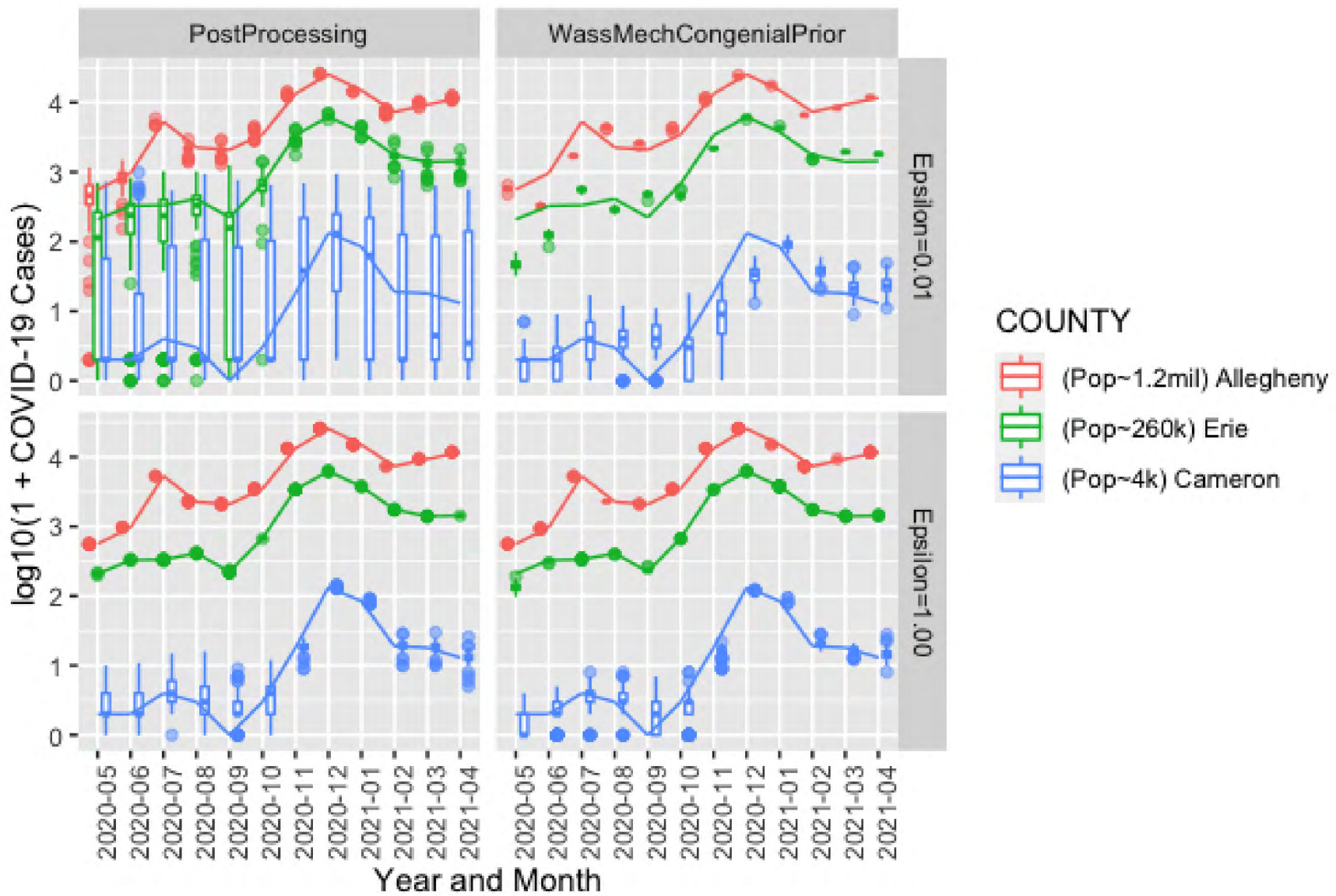
$$\nu_Z(Y) \propto \mathbb{1}_{\{Y \in \mathbb{Z}^J, \sum_{i=1}^J Y_{j,t}^{(c)} = s_t^c, Y_{j,t}^{(c)} \geq Y_{j,t}^{(d)} \geq 0\}}. \quad (39)$$

- 3 Prior congenial base measure:

$$\nu_Z(Y) \propto \phi(Y; s_t^{(c)}, x_{t-1}^{(c)}) \mathbb{1}_{\{Y \in \mathbb{Z}^J, \sum_{i=1}^J Y_{j,t}^{(c)} = s_t^c, Y_{j,t}^{(c)} \geq Y_{j,t}^{(d)} \geq 0\}}, \quad (40)$$

where  $\phi$  is the PMF of the Dirichlet-Multinomial distribution:

$$\phi(Y_t^{(c)}; s_t^{(c)}, \alpha x_{t-1}^{(c)}) = \frac{\Gamma(\alpha s_t^{(c)}) \Gamma(s_t^{(c)} + 1)}{\Gamma((\alpha + 1) s_t^{(c)})} \prod_{j=1}^J \frac{\Gamma(Y_{j,t}^{(c)} + \alpha X_{j,t-1}^{(c)})}{\Gamma(Y_{j,t}^{(c)}) \Gamma(\alpha X_{j,t-1}^{(c)} + 1)}. \quad (41)$$

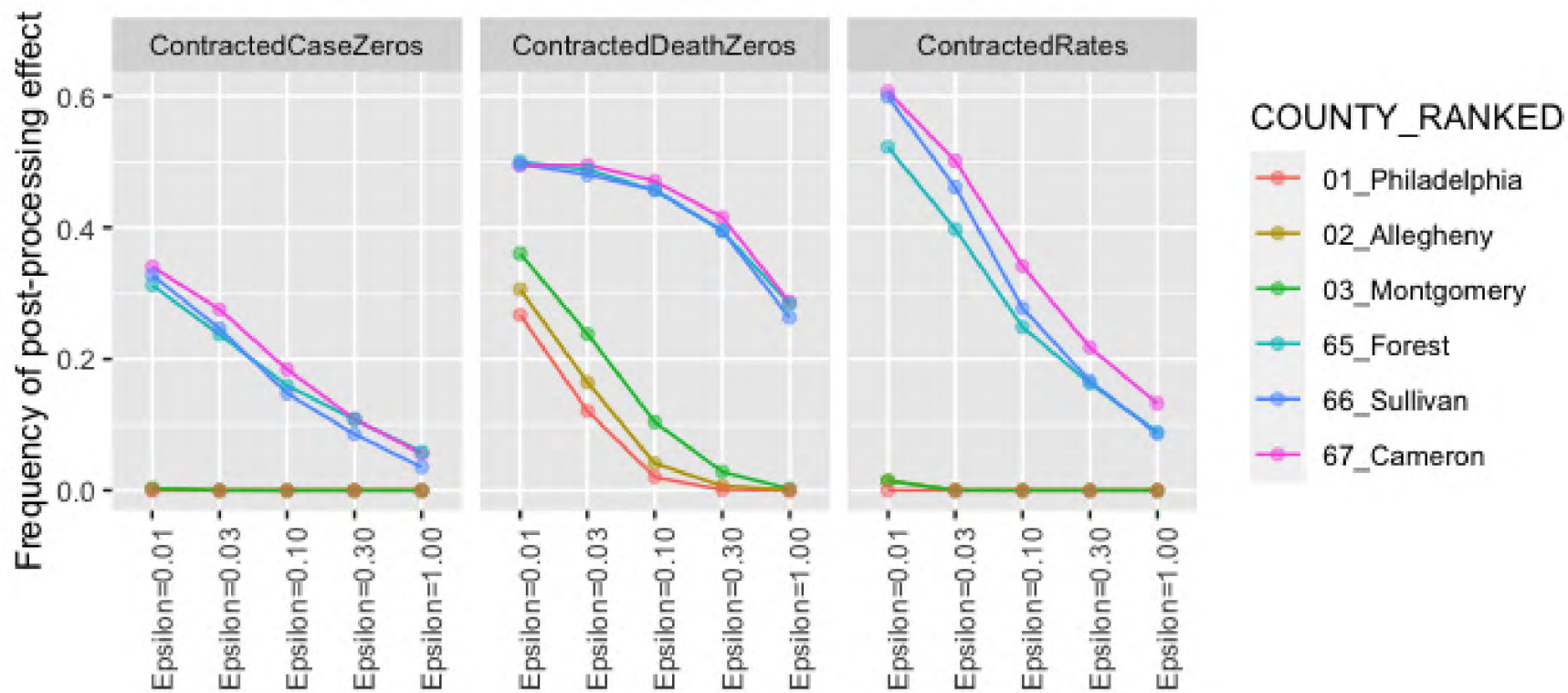


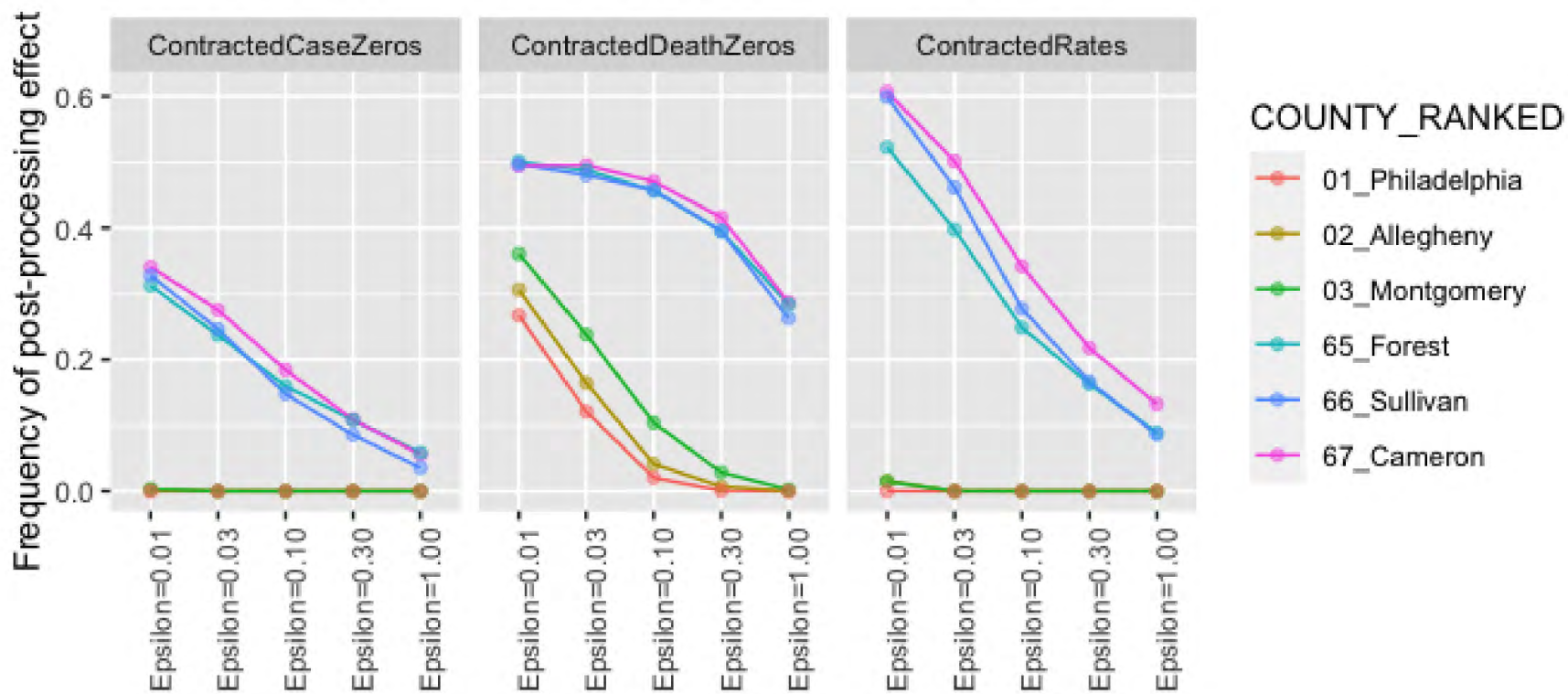
- In theory, Bayes estimators based on  $Y^*$  are worse than those based on  $Y, Z$ . However...
  - We can only numerically approximate the Bayes estimator based on  $Y^*$  due to intractable post-processing.
  - $\uparrow$  uncertainty quantification for estimators based on  $Y^*$  depend on the quality of this numerical approximation, potentially confounding the effects we want to isolate.
- Need an alternative strategy to compare estimators!

Theoretical problem: post-processing degrades statistical signal by mapping different sufficient statistics to the same output.

Empirical realizations for COVID-19 case study:

- `ContractedCaseZeros`: cases where multiple potential imputations of the private COVID-19 case data are contracted to 0
- `ContractedDeathZeros`: cases where multiple potential imputations of the private COVID-19 death data are contracted to 0
- `ContractedRates`: cases where the constraint that COVID-19 cases is bounded below by COVID-19 deaths contracts imputations of the COVID-19 survival rate to 0.





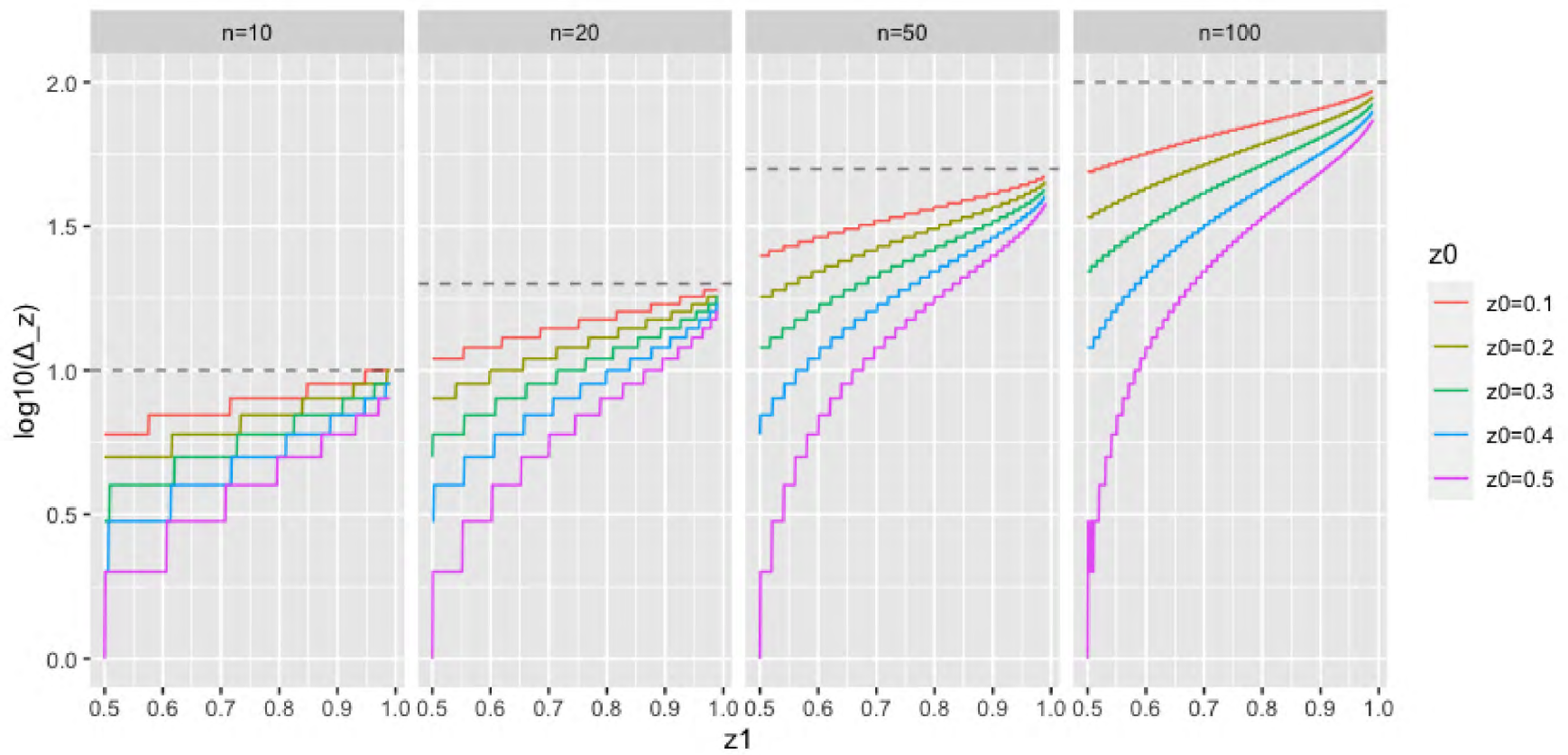
For binary count data, we consider families of distributions  $\Theta_{(z_0, z_1)}$ , indexed by two parameters  $z_1$  and  $z_0$ . Let  $i \neq j$  and assume  $Z$  defines the following probabilistic public information for all  $\theta_{(z_0, z_1)} \in \Theta_{(z_0, z_1)}$ :

$$\mathbb{P}(X_i = 1 \mid s_{j1}, Z = z) \leq z_1 \in [0, 1], \quad \mathbb{P}(X_i = 1 \mid s_{j0}, Z = z) \geq z_0 \in [0, 1]$$

Notes:

- As  $z_1 \rightarrow 1$  and  $z_0 \rightarrow 0$ , public information reveals more probabilistic information about private statistics.
- Related to the “worst-case” optimal transport solution, in which *all* the probability mass moves between two conditional distributions according to this dependency.





- 1 Discussion: Posterior-to-Posterior Semantics and PPD
- 2 Discussion:  $\epsilon$ -TP Robustness to Misspecification in  $\Theta_Z$
- 3 Discussion: MCMC Approximation Privacy Risks and Atomic Regeneration
- 4 Proof: WEM
- 5 Proof: WKNG
- 6 Proof: Exact Sampler Runtime
- 7 Proof: Exponential Family Stochastic Dominance
- 8 Algorithms: Perfect sampling
- 9 Algorithms: ABC Inference
- 10 Case study: Rural Alaska mortality
- 11 Case study: Pennsylvania Spatiotemporal COVID-19 data
- 12 Case study: worst-case MCMC convergence vs. realized exact runtime

## Key property

MCMC methods require accounting for the **slowest mixing** chain, but our method can be much faster because **the runtime depends on the realized confidential data**

Illustrative example: Laplace mechanism ( $L_X(y) = \|\bar{X} - y\|_1$ ) with data bounded in  $[0, 1]^d$

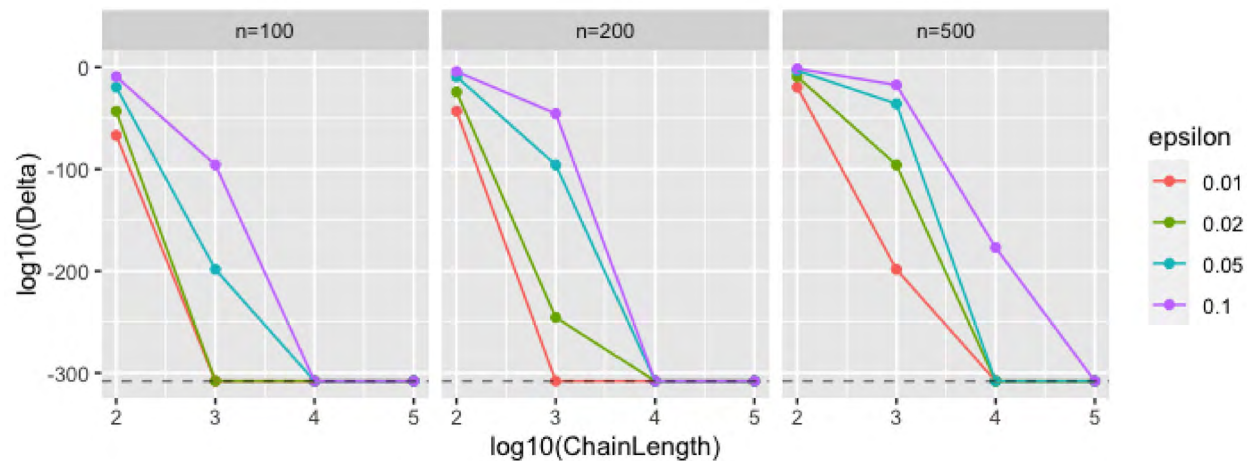
- Two original Markov chains: Metropolis-Hastings (MH) with independent uniform proposals and symmetric Laplace proposals with scale  $\alpha$
- Closed form expressions for worst-case  $\delta$  with MH MCMC (Mengersen and Tweedie, 1996)

$$\|\mu_{X,m} - \mu_X\|_{TV} \leq (1 - \beta_{\text{MCMC}})^m, \quad (42)$$

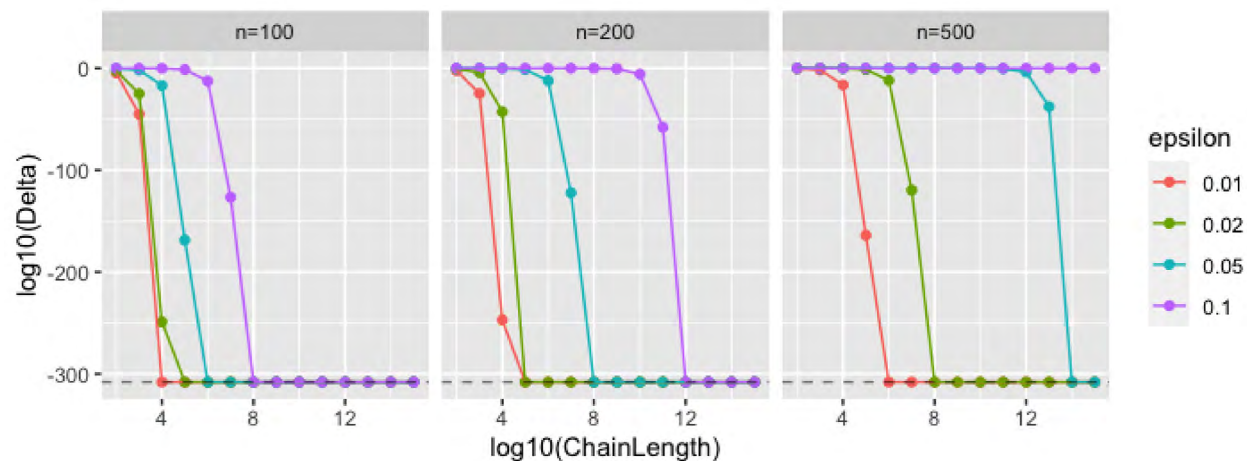
$$\begin{cases} \beta_{\text{MCMC,Unif}} \triangleq \left(\frac{2d}{\epsilon n} (1 - e^{-\epsilon n/2d})\right)^d \\ \beta_{\text{MCMC,Lap}} \triangleq (2\alpha)^d \exp\left(-\left(\alpha d + \frac{\epsilon n}{2}\right)\right) \left(\frac{1}{\alpha}(1 - e^{-\alpha})\right)^d \end{cases}$$

(Dashed line = 64-bit double precision threshold)

a) Independent uniform proposals:



b) Symmetric Laplace proposals, scale =  $\epsilon n/2$ :



## Benefits of our approach:

- Satisfies  $\epsilon$ -DP
- Runtime depends on realized confidential data, and *not* the confidential data for the slowest-mixing Markov Chain
- Only requires minorizing bound, and not properties of  $L_X$  (i.e. convexity, Lipschitz, etc.)
  - Most existing analyses of the exponential mechanism (like Ganesh and Talwar, 2020) require these assumptions
  - Demonstrates why regeneration is more suited to this problem than other perfect sampling methods like coupling from the past (Propp and Wilson, 1996)
- Easily extendable to other MC algorithms satisfying minorization condition

## Limitations of our approach:

- Uniform ergodicity assumption: methods do not have finite expected runtime for unbounded state spaces, like  $\mathbb{R}^d$ . Caveats:
  - $\mathcal{Y}$  is often artificially restricted to bound the sensitivity of  $L_X$ , so this issue is not prevalent in practice.
  - For some perfect sampling algorithms (like ours based on Lee et al, 2014), uniform ergodicity is a **necessary** requirement for finite expected runtime.
- Minorizing constant suffers from curse of dimensionality
- Side-channel vulnerability: multiple replications of similar queries could leak information about confidential data through runtime

Traditional analysis: privacy vs. utility

## Extensions of our work: privacy vs. utility vs. runtime

- Trading off utility and runtime:
  - Exponential mechanisms can be implemented exactly over enumerable discrete state spaces
  - $\implies$  corollary: if we release a sample from a discrete approximation w.p.  $k$ , then we reduce runtime at the cost of some utility
- Trading off privacy and runtime:
  - (Awan and Rao, 2021) consider rejection sampling where  $N_{\text{prop}}$  is known and can leak information
  - $\implies$  corollary: with longer artificial runtime, can release  $\tilde{N}_{\text{prop}} \perp\!\!\!\perp X$  with 0-DP so that  $(Y, \tilde{N}_{\text{prop}})$  is  $\epsilon$ -DP

Setup:  $\mathcal{Y} \triangleq \{y \in \mathbb{R}^p \mid \|y\|_1 \leq B\}$ ,  $X \in [-1, 1]^{n \times p}$ ,  $Z \in [-1, 1]^n$ . Define the original loss function  $\ell_X(y)$ :

$$\ell_X(y) = \frac{1}{2} \|Z - Xy\|_2^2 + \frac{\lambda}{2} \|y\|_2^2.$$

For any two  $X, X'$  that differ on one element and all  $y \in \mathcal{Y}$ :

$$\|\nabla \ell_X(y) - \nabla \ell_{X'}(y)\|_1 \leq 2(1 + B) \sup_{x \in [0, 1]^p} \|x\|_1 \leq 2(1 + B)p$$

This yields a final mechanism given by:

$$f_X(y) \propto \exp\left(-\frac{\epsilon}{4(1 + B)p} \|X^T(Xy - Z) + \lambda y\|_1\right) \mathbb{1}_{\{\|y\|_1 \leq B\}}. \quad (43)$$



Runtime relaxation: sample from  $Y^* \sim \tilde{\mu}_X$ , and if  $Y^* = a$ , sample  $Y \sim \mu_X^{(\text{disc})}$

- $\mu_X^{(\text{disc})}$  implemented over  $\ell$  sample points iid from  $\text{Unif}(\mathcal{Y})$
- Reduces runtime by preventing multiple draws from  $\tilde{\mu}_X$

How to measure loss in utility?

$$\text{Err}(\epsilon, \ell, \varepsilon) \triangleq Q_{.05} \{ \mathbb{P}(L_X(y)^{(n_{\text{exp}}, \epsilon, \ell)} \geq \varepsilon) \}.$$

where  $Q_{.05}$  is the 5% empirical quantile across  $n_{\text{exp}}$  experiment replications.

Constants:

$$\begin{cases} n \triangleq 100 \\ p \triangleq 5 \\ \beta \triangleq (.1, .2, -.3, 0, 0)^T \\ \lambda \triangleq 1 \end{cases}$$

Random variables:

$$\begin{cases} X_{ij} \sim \text{Beta}(5, 5) & i \in [n], j \in [p] \\ e_i \sim \text{Beta}(20, 20) & i \in [n] \\ Z_i \triangleq X_{i,\cdot} \beta + (2e_i - 1) & i \in [n] \end{cases}$$

- Moderately-sized discrete approximations of the state space can provide comparable utility while saving runtime.
- Relative effect depends on privacy budget  $\epsilon$  and error tolerance  $\varepsilon$
- Still suffers from curse of dimensionality ( $n = 100, p = 5$ )

