

## Semi Differential Privacy

Young Hyun Cho and Jordan Awan  
Department of Statistics, Purdue University

### Motivation

#### Differential Privacy and the 2020 Census

How will differential privacy affect 2020 census data?

-Differential privacy has strong guarantees, but is not always implemented exactly

-For example, US Census publishes a combination of DP and exact statistics

-Currently NO satisfactory method of quantifying the privacy in such settings

### Preliminary

A randomized mechanism  $M(X)$  is  $(\epsilon, \delta)$ -Differential Privacy if for any  $\epsilon, \delta > 0$ ,

$$\mathbb{P}(M(X) \in S) \leq e^\epsilon \mathbb{P}(M(X') \in S) + \delta,$$

$\rho$ -zCDP if for any  $\rho \geq 0$  and for all  $\alpha \in (1, \infty)$ ,

$$D_\alpha(M(X) || M(X')) \leq \rho \alpha,$$

where  $D_\alpha$  is  $\alpha$ -Renyi divergence, for any databases  $X, X'$  such that  $d(X, X') \leq 1$

Smaller Privacy Parameters  $(\epsilon, \delta)$  or  $\rho \Rightarrow$  Stronger Privacy

BUT THE SAME SPIRIT:

Any Pair of Input Data  $X, X'$  with ONE DIFFERENT ENTRY is Difficult to distinguish by observing the output of  $M$ .

-What happens if true statistic  $T(X)$  is published along with DP  $M(X)$ ?  
-Especially, what if there is no neighboring data  $X, X'$  s.t.  $T(X) = T(X')$ ? That is,  
 $\{(X, X') : T(X) = T(X') \text{ and } d(X, X') \leq 1\} = \emptyset$

We propose a framework semi-DP, which properly accounts for the combination of private and non-private releases.  
Under semi-DP we derive optimal mechanisms under a variety of scenarios.

### Definition of Semi-DP

Algorithm 1 Semi-adjacent Parameter  $a$

- 1: **Input:** True Statistics  $T(X)$
- 2: Define  $\mathcal{X}^n|_T = \{X' \in \mathcal{X}^n : T(X') = T(X)\}$
- 3: **for**  $i \in [n]$  **do**
- 4: Define  $\mathcal{P}'_i = \{x'_i \in \mathcal{X} : \exists Y \in \mathcal{X}^{n-1} \text{ s.t. } T\left(\begin{bmatrix} x'_i \\ Y \end{bmatrix}\right) = T\left(\begin{bmatrix} x_i \\ X_{-i} \end{bmatrix}\right)\}$
- 5: **if**  $\mathcal{P}'_i = \emptyset$  **then**
- 6: Put  $a_i = \infty$
- 7: **else**
- 8: Define  $\mathcal{X}'_i = \left\{ \left( \begin{bmatrix} x'_i \\ Y \end{bmatrix} \right) : x'_i \in \mathcal{P}'_i, Y = \operatorname{argmin}_{Y \in \mathcal{X}^{n-1}} d(X_{-i}, Y) \text{ s.t. } T\left(\begin{bmatrix} x'_i \\ Y \end{bmatrix}\right) = T(X) \right\}$
- 9:  $a_i = \max_{X'_i \in \mathcal{X}'_i} d(X, X'_i)$
- 10: **end if**
- 11: **end for**
- 12: **Output:**  $a = \max_{i \in [n]} a_i$

Given semi-adjacent parameter  $a$ , our semi-DP makes difficult to distinguish any pair of data in

$$\{(X, X') : T(X) = T(X') \text{ and } d(X, X') \leq a\}$$

### Analysis on US 2020 Census

$T(X)$ : True Population for each state, the District of Columbia and Puerto Rico

$\Rightarrow$  semi-adjacent parameter  $a = 2$

Census Advertised:  $\rho$ -zCDP with  $\rho = 0.213$

Semi-DP framework:  $\rho$ -zCDP with  $\rho = 0.852$

### Mechanism Design

Goal: Add the least noise and achieve the target privacy parameter

Algorithm 2 Generalized Optimal  $K$ -norm mechanism:  $\epsilon$ -semi DP

- 1: **Input:** Statistics  $T(X)$  to be released, Mechanism  $\theta$  to be privatized
- 2: Calculate the semi-adjacent parameter  $a$
- 3: Define the sensitivity space  $S_{\theta|T} = \{\theta(X) - \theta(X') : d(X, X') \leq a, T(X) = T(X')\}$
- 4: Define  $\operatorname{Proj}_{S_{\theta|T}}$  be the orthogonal projection operator on to  $\operatorname{span}(S_{\theta|T})$
- 5: Calculate  $S_{\operatorname{Proj}} = \{\operatorname{Proj}_{S_{\theta|T}}\theta(X) - \operatorname{Proj}_{S_{\theta|T}}\theta(X') : d(X, X') \leq a, T(X) = T(X')\}$
- 6: Let  $\operatorname{Hull}(S_{\operatorname{Proj}})$  be a convex hull of  $S_{\operatorname{Proj}}$  and define a norm  $\|\cdot\|_K$  by  $\|u\|_K = \inf\{c \in \mathbb{R}^{\geq 0} : u \in cH\}$
- 7: Calculate  $K$ -norm sensitivity  $\Delta_K = \sup_{u \in S_{\operatorname{Proj}}} \|u\|_K$
- 8: Sample  $V \sim f_V(v) = c \exp\left(-\frac{\epsilon}{\Delta_K} \|v\|_K\right)$
- 9: **Output:**  $M(X) = \theta(X) + V$

Generalized Optimal  $K$ -norm mechanism achieves  $\epsilon$ -semi DP with the least additive noise, so optimal!

Algorithm 3 Gaussian Mechanism:  $\mu$ -semi GDP

- 1: **Input:** Statistics  $T(X)$  to be released, Mechanism  $\theta$  to be privatized
- 2: Calculate the semi-adjacent parameter  $a$
- 3: Define the sensitivity space  $S_{\theta|T} = \{\theta(X) - \theta(X') : d(X, X') \leq a, T(X) = T(X')\}$
- 4: Calculate  $P$  be the orthogonal projection matrix onto  $\operatorname{span}(S_{\theta|T})$
- 5: Calculate  $l_2$ -norm sensitivity  $\Delta_2 = \sup_{u \in S_{\theta|T}} \|u\|_K$
- 6: Sample  $N \sim N(0, (\frac{\Delta_2}{\mu})^2 P)$
- 7: **Output:**  $M(X) = \theta(X) + N$

Gaussian Mechanism achieves  $\mu$ -semi GDP with less noise than naïve approaches

Additive Gaussian noise can also achieve  $(\epsilon, \delta)$ -DP and  $\rho$ -zCDP