

Today

• Review: inverse sensitivity

• Privacy Wrapper Theorem

• Monotone Functions & Interior Points

Lecture 25

4-21-26

Inverse Sensitivity Mechanism

Goal: approximate function  $f: \mathcal{X}^n \rightarrow \mathcal{Y}$ ,  $|\mathcal{Y}| < \infty$

$$P(y) \propto \exp\left(-\frac{\varepsilon}{2} \text{len}_f(y; x)\right)$$

$$\text{len}_f(y; x) = \min_{x'} \{d_H(x, x') \mid f(x') = y\}$$

Claim: This is  $\varepsilon$ -DP.

Proof Fix  $y$ ,  $x \sim x'$ . Then

$$\text{len}_f(y; x) = \min_{x''} \{d_H(x, x'') \mid f(x'') = y\}$$

$$\leq \min_{x''} \{d_H(x', x'') + d(x, x') \mid f(x'') = y\}$$

$$\leq 1 + \text{len}_f(y; x')$$

Sensitivity bounded by one  $\Rightarrow$  exponential mechanism is  $\epsilon$ -DP  $\square$

Canonical application: median

Accuracy guarantee:

Def The  $k$ -neighborhood of  $x$  is

$$\mathcal{N}_k(x) = \{x' : d_H(x, x') \leq k\}$$

Theorem [Asi & Duchi 2020] For  $k = \frac{2}{\epsilon} \log\left(\frac{2|\mathcal{X}|^2}{\beta\epsilon}\right)$ , w.p.  $1 - \beta$  the

inverse sensitivity mechanism returns  $\tilde{y}$  such that

$$\min_{x' \in \mathcal{N}_k(x)} f(x') \leq \tilde{y} \leq \max_{x' \in \mathcal{N}_k(x)} f(x)$$

Very simple algorithm!

Optimal in a strong sense, but

major drawbacks:

- $\square$  unclear how to compute  $\text{len}_f(y|x)$  when  $|\mathcal{X}| = \infty$
- $\square$  for many tasks, asymmetry between removal & addition; when you

can add points, can add outliers.

## Privacy Wrappers

Def The k-down neighborhood of  $x$  is

$$\mathcal{N}_k^\downarrow(x) = \{x' : x' \leq x, |x'| \geq |x| - k\}$$

Theorem [Linder, Raskhodnikova, Smith, Steinke 2025]

There is an  $(\epsilon, \delta)$ -DP mechanism that,

for any  $f: \mathcal{X}^n \rightarrow \mathcal{Y}$ , for

$$k = \frac{\log(1/\delta)}{\epsilon} \cdot 2^{O(\log^4(1/\delta))},$$

returns  $\tilde{y}$  such that

$$\min_{x' \in \mathcal{N}_k^\downarrow(x)} f(x) \leq \tilde{y} \leq \max_{x' \in \mathcal{N}_k^\downarrow(x)} f(x')$$

Moreover, this algorithm only queries  $f$  on datasets in  $\mathcal{N}_k^\downarrow(x)$ ; running time

$$n^{O(\frac{1}{\epsilon} \log(1/\delta))}. \text{ (time to compute } f)$$

Remarks:

□ Algorithm more complicated than inverse sensitivity

□ Fixes drawbacks: always computable with query access to  $f$ , always accuracy no worse.

□ LRSS25 also give pure-DP version which succeeds w.p.  $1-\beta$  and has  $k = O(\frac{1}{\epsilon} \log(1/\beta))$

□ Looks familiar to algorithmic robust statistics: "for all large subset..."

## Open Directions

- What is the analog of this for multivariate outputs?
- When can this be implemented efficiently?
- When can we accommodate  $|Y| = \infty$ ?

## Analysis: Monotonicity & Shifted Inverse

Fang, Dong, Yi 2022 give an algorithm for monotone functions

Def A function  $f: X^n \rightarrow Y$  is monotone if  $x \leq x'$  implies  $f(x) \leq f(x')$ .

Theorem [FDY 22] For any monotone  $f: \mathcal{X}^n \rightarrow \mathcal{Y}$ ,  
 the Shifted Inverse (Sensitivity) Mechanism  
 is  $\epsilon$ -DP and for  $k = O(\frac{1}{\epsilon} \log \frac{|\mathcal{Y}|}{\rho})$  w.p.  $1 - \beta$   
 returns  $\tilde{y}$  such that

$$\min_{x' \in \mathcal{N}_k^\downarrow(x)} f(x') \leq \tilde{y} \leq f(x)$$

Shifted inverse defines loss

$$l_f(y; x) = \min \{ d_H(x, x') : x' \leq x, f(x') \leq y \}$$

Claim For monotone  $f$ , any  $y \in \mathcal{Y}$ , and  $x \sim x'$ ,

$$|l_f(y; x) - l_f(y; x')| \leq 1$$

Proof: elegant, go look it up!  $\square$

Run exponential mechanism to find  $\tilde{y}$   
 approximately minimizing

$$|l_f(\tilde{y}; x) - k|$$

For the main result in LRSS25, take this algorithm and

1) median  $\Rightarrow$  "generalized interior point"

2) Take any  $f$  and "monotonize" it.

Informally:

$$M_\ell[f](x) = \max \{ f(s) : s \leq x, |s| \geq |x| - \ell \}$$

---

More detail on inverse sensitivities

---

The exponential mechanism we've seen for the median uses loss

$$\ell_{\text{median}}(y; x) = \left| \frac{n}{2} - \#\{i : x_i \leq y\} \right|$$

This is exactly "the number of points we need to change in  $x$  to make  $y$  a median point." Thus, we can see the connection to inverse sensitivity.

Here is another view on the median's loss function: it's a shifted version of the maximum. We have

$$l_{\max}(y; x) = \#\{i : x_i \leq y\}$$

as a measure of how well  $y$  approximates the maximum of  $x$ . We measure the approximation by counting data points.

Note that the maximum is a monotone function. We can regard the median\* as a shifted version of the maximum. \*or any quantile

FDY use loss function

$$l_f(y; x) = \min \{d(x, x') : x' \in x, f(x') \leq y\}$$

and run the exponential mechanism to minimize  $(l_f(y; x) - k)$  for

some parameter  $k$ .

We've already seen this! Our loss function for the max is

$$l_{\max}(y; x) = \#\{i : x_i \leq y\}$$

$$= \min \{d(x, s) : s \leq x, \max(s) \leq y\}$$

And we've seen how to shift it.

## References

Asi & Duchi 2020

<https://arxiv.org/abs/2005.10630>

Fang, Dong, & Yi 2022

<https://juanru-fang.github.io/ShiftedInverse.pdf>

Steinke 2023 blog post

<https://differentialprivacy.org/down-sensitivity/>

LRSS 2025

<https://arxiv.org/abs/2503.19268>

