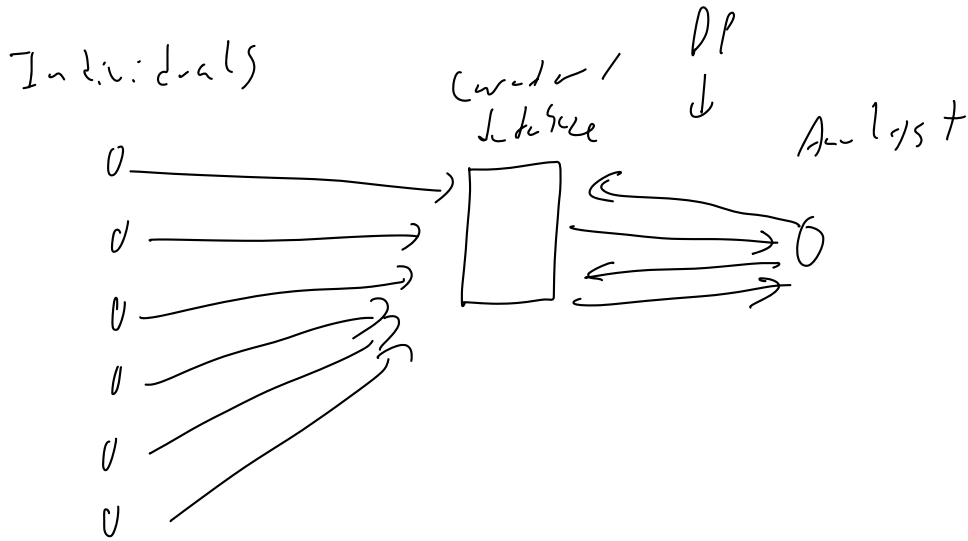


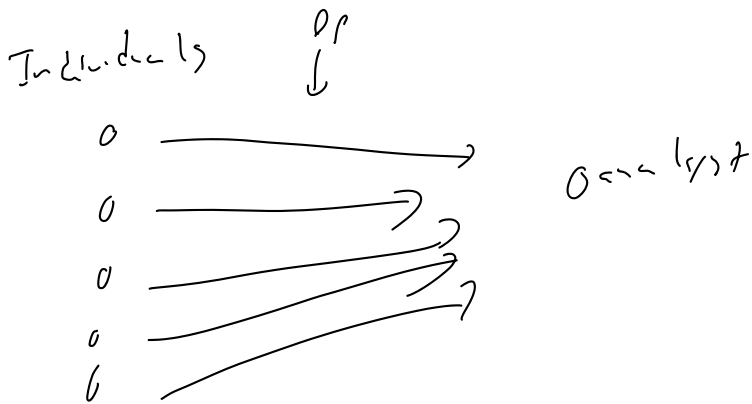
Local DP:

so far: centralized DP. "Trusted curator?"



Analyst interacts w/ curator to produce privacy, but  
curator knows all true values!  
- what if don't trust curator to have your data?  
- what if curator gets hacked?

Local DP: don't even trust curator!



Formal def: "local randomizers", "local random coins", etc.

completely equivalent to every individual using DP on dataset of size 1!

Def: If every individual uses an  $(\epsilon, \delta)$ -DP alg,  
then combined alg/mechanism/protocol is  $(\epsilon, \delta)$ -local DP

Thm:  $(\epsilon, \delta)$ -local DP  $\Rightarrow$   $(\epsilon, \delta)$ -DP

Pf: Parallel composition!

Variants: - interactive vs. non-interactive

Pros: - super distributed, no trust, - used in practice!

Cons: - often much worse accuracy!

Easy example: sum of binary values

- Each  $x_i \in \{0, 1\}$ , return  $X = \sum_{i=1}^n x_i$

- Central DP:  $\Delta f = 1 \Rightarrow \tilde{X} = X + \text{Lap}(\sqrt{\epsilon})$  :  $\epsilon$ -DP

$$\mathbb{E}[|X - \tilde{X}|] = O(\sqrt{\epsilon})$$

- Local DP: only possible algorithm is randomized response!

- How: construct  $\tilde{x}_i = \begin{cases} x_i & \text{w/prob. } \frac{e^\epsilon}{e^\epsilon + 1} \\ 1-x_i & \text{w/prob. } \frac{1}{e^\epsilon + 1} \end{cases}$

Claim:  $\epsilon$ -local DP.

PF: (unfair) under i.

$$\frac{\Pr[\tilde{x}_i = 1 \mid x_i = 1]}{\Pr[\tilde{x}_i = 1 \mid x_i = 0]} = \frac{\frac{e^\tau}{e^\tau + 1}}{\frac{1}{e^\tau + 1}} = e^\tau$$

$$\frac{\Pr[\tilde{x}_i = 0 \mid x_i = 0]}{\Pr[\tilde{x}_i = 0 \mid x_i = 1]} = e^\tau \quad \checkmark$$

Given  $\tilde{x}_i$ 's, what should (unbiased) estimator do?

$$E\left[\sum_{i=1}^n \tilde{x}_i\right] = \sum_{i: x_i=1} \frac{e^\tau}{e^\tau + 1} + \sum_{i: x_i=0} \frac{1}{e^\tau + 1}$$

$$\begin{aligned} \Rightarrow (e^\tau + 1) E\left[\sum_{i=1}^n \tilde{x}_i\right] &= \sum_{i: x_i=1} e^\tau + \sum_{i: x_i=0} 1 \\ &= \sum_{i=1}^n 1 + \sum_{i: x_i=1} (e^\tau - 1) \\ &= n + \sum_{i: x_i=1} (e^\tau - 1) \end{aligned}$$

$$\Rightarrow \frac{(e^\tau + 1) E\left[\sum_{i=1}^n \tilde{x}_i\right] - n}{e^\tau - 1} = \sum_{i: x_i=1} 1$$

$$\Rightarrow E\left[\sum_{i=1}^n \left(\frac{e^\tau + 1}{e^\tau - 1} \tilde{x}_i - \frac{1}{e^\tau - 1}\right)\right] = \sum_{i: x_i=1} 1 = \#1's.$$

So, estimator (unbiased)

How good is this?

$$\text{Let } X \sim \#1's, \quad \tilde{x} = \sum_{i=1}^n \left( \frac{e^{\tilde{x}_i \gamma}}{e^{\tilde{x}_i \gamma} + 1} \tilde{x}_i - \frac{1}{e^{\tilde{x}_i \gamma} + 1} \right)$$

Variance:  $E[(X - \tilde{x})^2] \approx \frac{n}{\gamma^2} \quad (\text{when } \gamma \text{ small}) \approx \sigma^2$

$$\Rightarrow E[|X - \tilde{x}|] \geq \frac{\sqrt{n}}{\gamma}, \text{ which is in hypothesis/choice}$$

So constant expected loss in control,  $\sqrt{n}$  in local!

Private Selection:

$$\text{Each } x_i \in (d), \text{ return } \arg\max_{j \in F(d)} \left( \sum_{i: x_i = j} 1 \right)$$

Control: RNM: add  $(\tilde{x}_i, (1/\gamma))$  to each entry return max.

$$\text{threshold, if diff b/w next, second} \geq \frac{\log d}{\gamma}$$

Localize: write each elt of  $(d)$  as  $\{c_i\}^{\log d}$ , then

RR on each of  $\log n$  bits.

Composition: let  $\epsilon = \frac{\gamma}{\log d}$

$$\Rightarrow \text{error of each bit} \approx \frac{\sqrt{n} \log d}{\gamma}$$

$$\Rightarrow \text{need } \underline{\text{every bit}} \text{ to have } \geq \frac{\sqrt{n} \log d}{\gamma} \text{ more of correct value.}$$

Could be pretty unlikely even if lots more of correct next!

Ex:  $d=4$   $\begin{matrix} 00 \\ 01 \\ 10 \\ 11 \end{matrix}$   $\begin{matrix} 4 \\ 10 \end{matrix}$   $\begin{matrix} 10 \\ 01 \end{matrix}$   $\begin{matrix} 11 \\ 11 \end{matrix}$

$\Rightarrow$  final bit will be set to 1 incorrectly!

- "One-hot" encoding: each bit of  $(d)$  written as  $\{c, 1\}^d$   
 $x \in d \rightarrow \{c, 0, \dots, c, \underset{\substack{\uparrow \\ x}}{1}, c, \dots, c\}$

RR on each bit: don't need comparison! Like RNN

$\Rightarrow$  expected error of each bit  $\sim \frac{\sqrt{n}}{\epsilon}$

$\Rightarrow$  ok if diff b/w largest, second largest  $\geq \frac{\sqrt{n} \log d}{\epsilon}$

More thoughts:

- Corrections to "statistical query alg's"
- Used significantly in industry!
- Google: RAPPOR
  - "Permanence" RR
  - Bloom filter to map each  $x \in X \rightarrow \{c, 1\}^k$
  - Issue: assumes value never changes!
- Apple: has a "count sketch" alg. to reduce comm. cost
- Microsoft:
  - more tricks to reduce comm. cost
  - What if values change over time, but not by much?

- e.g. how long spent in app for day
- discretize!
- too fine: too much re-randomization
- too coarse: low utility
- Main result: 3-bit mean estimation w/ correctness, as real loss in accuracy!
- In graphs: LEOP pretty powerful!
- one edge changes  $\Rightarrow$  two nodes have diff. local input
- but not arbitrarily different!
- Interactivity super helpful.
- Not just PR!

### Beyond Local:

- Crypside: replace trusted curator by secure multiparty computation! (MPC)
- basically works, but need all users to be available to run protocol
- computational assumption, PKC (public key infra)

Next time: "shuffle model".