

Binary Tree Mechanism

- ◊ Problem formulation: gallery release
- ◊ Algorithm: How to add "smart" noise?
- ◊ Applications: quantiles, CDF, continual release.

Algorithms So far ...

Answering Numeric Queries

- Laplace Mechanism ↗^{more general}
 ↖^{less noise}
- Randomized Response



Focus.

Selection Problem ↵

- Exponential Mechanism ↵
- Report Noisy Max. ↵

Query Release Problem.

Example: Census Bureau

Focus: Linear / counting queries

Given $x = (x_1, \dots, x_n) \in X^n$

$$f(x) = \sum_{i=1}^n \varphi(x_i) \quad \text{for } \varphi: X \mapsto \{0, 1\}$$

\uparrow
 Specify some property
 e.g., (Asian, 30's, male)

Goal: Release answers $\{a_1, \dots, a_k\}$ for queries $\{f_1, \dots, f_k\}$
 with accuracy δ & ϵ -DP.

$$\max_{j=1}^k |a_j - f_j(x)| \leq \alpha n$$

\uparrow
 e.g., 1%

Applying Laplace Mechanism.

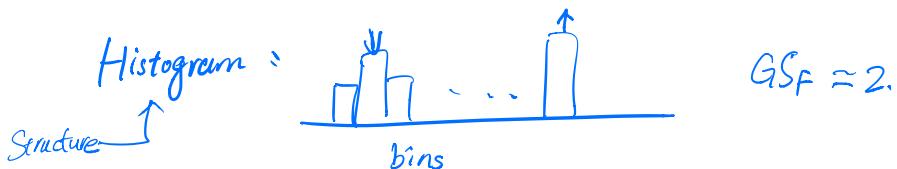
Define $F : X^n \mapsto \mathbb{R}^k$ change some κ_i
 s.t. $F(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_k(x) \end{pmatrix} \leftarrow \pm 1$

ℓ_1 -sensitivity $GS_F := \max_{x \sim x'} \|F(x) - F(x')\|_1$

$$\boxed{GS_F = k \quad \text{in worst case.}}$$

Lap Mech: $f_j(x) + Z_j$, $Z_j = \text{Lap}\left(\frac{k}{\epsilon}\right)$

max error $\leq O\left(\frac{k \log k}{\epsilon}\right) \leq \alpha n$.



Interval / Threshold Queries.

Given dataset $x = (x_1, \dots, x_n) \in X$

Data Universe
"Domain" $X = \{1, \dots, D\} \equiv [D]$

↳ e.g., height

time arrival

salary. ---

($D \gg n$)

Interval Queries : given $1 \leq s \leq t \leq D$

$$f_{s,t}(x) = \# \{i \mid s \leq x_i \leq t\}$$

Threshold Queries : given $t \in [D]$

$$f_t(x) = \# \{i \mid x_i \leq t\}$$

Examples :

① CDF (cumulative distribution function) : $\Phi(t) = \# \{i \mid x_i \leq t\} / n$

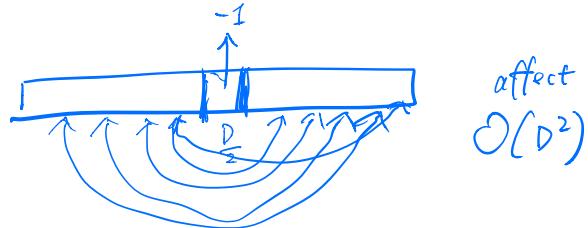
② g -quantile : threshold t such that $\Phi(t) = g$. ($\Phi^{-1}(g)$)

Baseline Laplace Mechanism

Answering all $\binom{D}{2} = O(D^2)$ interval queries

$$F(x) = \begin{pmatrix} f_{s,t}(x) \\ \vdots \\ f_{s,D}(x) \end{pmatrix}_{1 \leq s \leq t \leq D}$$

what is GS_F ?

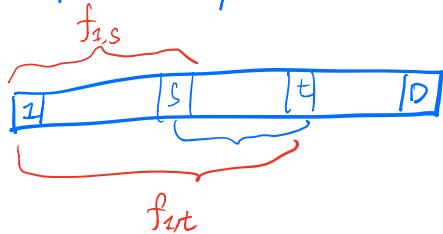


$\forall (s,t) :$

$$\alpha_{s,t} = f_{s,t}(x) + Z_{s,t}, \quad Z_{s,t} \sim \text{Lap}\left(\frac{GS_F}{\epsilon}\right) \in O\left(\frac{D^2}{\epsilon}\right)$$

$$\mathbb{E}\left[\max_{(s,t)} |\alpha_{s,t} - f_{s,t}(x)|\right] = O\left(\frac{D^2 \log D}{\epsilon}\right) \leftarrow \text{2n error.}$$

Simple Improvement over Basic Laplace.



$$f_{s,t} = f_{1,t} - f_{1,s}$$

$\leq \leq$
 $a_{1t} \quad a_{1s}$

$\forall 1 \leq t \leq D$

$$f_{1,t} + Z_{1,t}, \quad Z_{1,t} \sim \text{Lap}\left(\frac{D}{\epsilon}\right)$$

$$\mathbb{E}[\text{max-error}] \leq O\left(\frac{D \log D}{\epsilon}\right) \quad \leftarrow \begin{matrix} \text{previously} \\ O\left(\frac{D^2 \log D}{\epsilon}\right) \end{matrix}$$

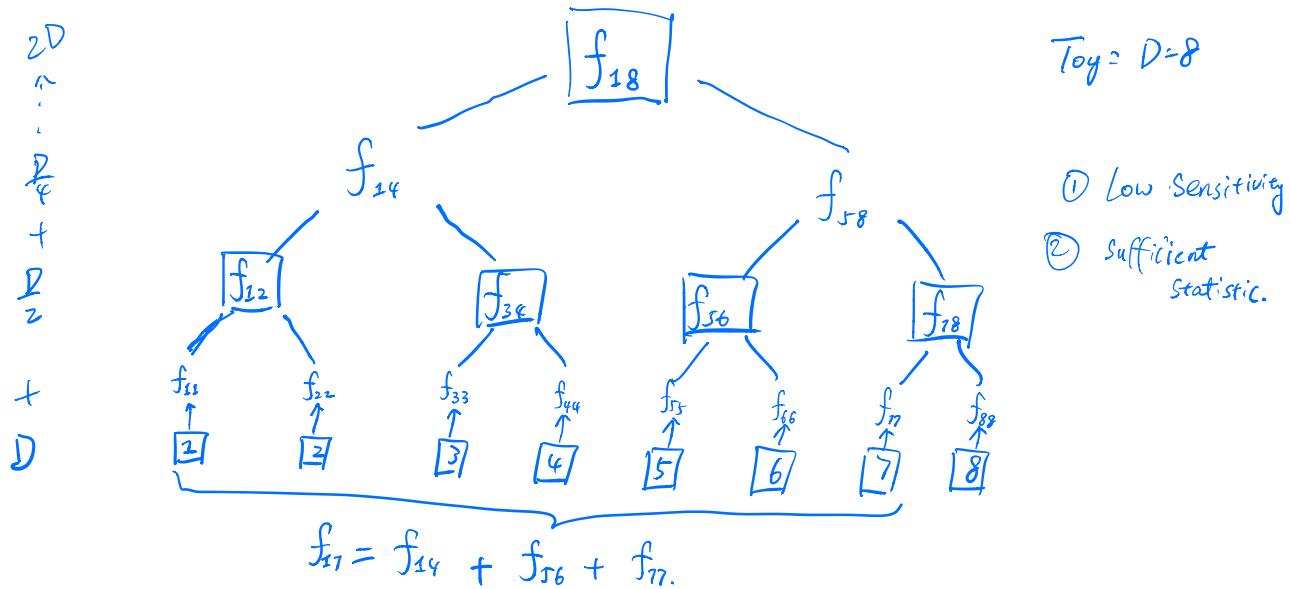
Answer an easier set of sufficient statistics.

→ post-process to get answers.

Binary Tree Mechanism

$$\text{error} \rightarrow \frac{1}{\epsilon} \log^3 D$$

$$Wlog \Rightarrow D = 2^m$$



$$T = \left\{ (s, t) : s = j \cdot 2^{l-1} + 1 \text{ and } t = (j+1) \cdot 2^{l-1} \right\}$$

for $1 \leq l \leq \log_2 D$ and $1 \leq j \leq D/2^{l-1}$.

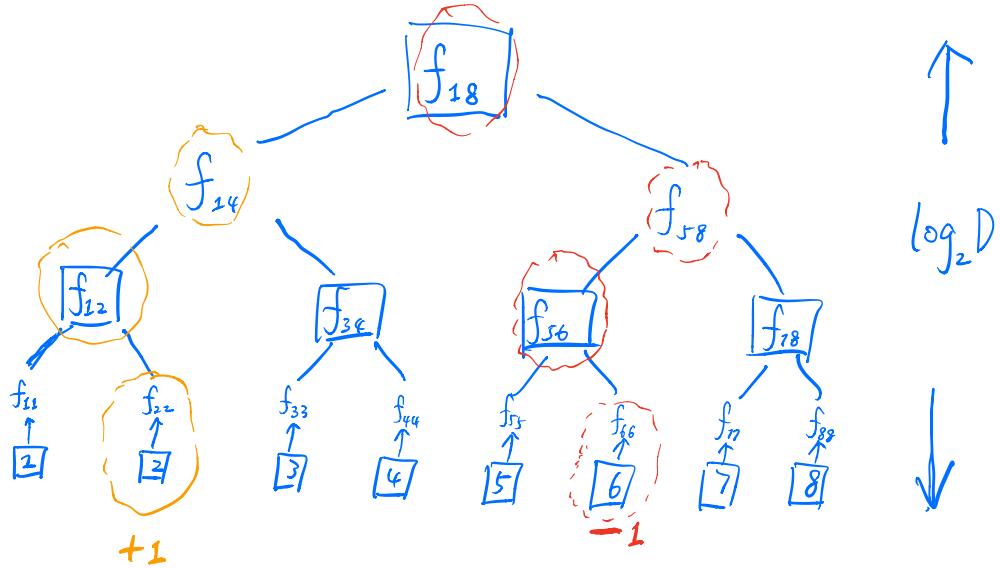
What is the sensitivity F_T ? $\approx \log_2 D$

$$|T| = \sum_{l=1}^{\log_2 D} \frac{D}{2^{l-1}} = 2D - 1.$$

Claim : Let $F_T^{(x)} = (f_{s,t}(x))_{(s,t) \in T}$

GS_F is at most $2^{\log_2 D}$.

Proof.



Binary Tree Mechanism

Release tree node intervals.

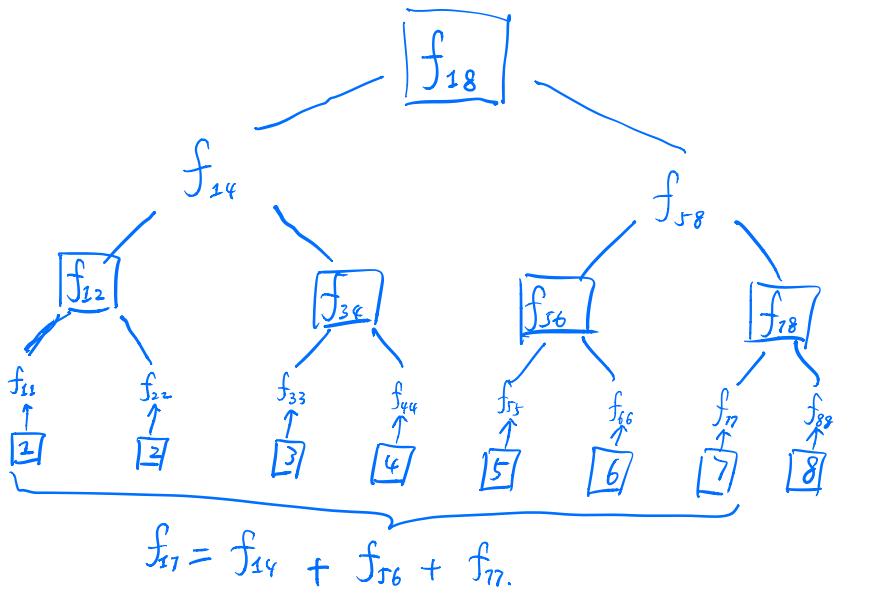
for every $(s,t) \in T$ do

$$a_{s,t} = f_{s,t}(x) + z_{s,t},$$

$$z_{s,t} \sim \text{Lap}\left(\frac{2 \log D}{\epsilon}\right) \in O\left(\frac{\log D}{\epsilon}\right)$$

$$\mathbb{E}\left[\max_{(s,t) \in T} |z_{s,t}|\right] = O\left(\frac{\log^2(D)}{\epsilon}\right).$$

Recovering the Answers.



Claim. For every $1 \leq t \leq D$, there exists $S \subseteq T$ of size $|S| \leq \log_2 D$ such that

$$f_{t,t} = \sum_{(u,v) \in S} f_{u,v} \rightarrow \approx \sum_{u,v}$$

Privacy + Accuracy Guarantees.

Theorem. Bin Tree Mechanism is $\boxed{\epsilon\text{-DP}}$

and releases $(a_{s,t})_{1 \leq s \leq t \leq D}$ such that

$$\mathbb{E} \left[\max_{(s,t)} |a_{s,t} - f_{s,t}(x)| \right] \leq O\left(\frac{1}{\epsilon} \log^3 D\right).$$

Proof Sketch = $\epsilon\text{-DP} \Leftarrow \text{Lap Mech w/ } \frac{2 \log(D)}{\epsilon}$ sensitivity

Accuracy: ① Suffices to release $f_{s,t} \forall t \in [D]$



$$② f_{1,t} = \sum_{(u,v) \in S} f_{u,v} \leftarrow$$

$$\mathbb{E} \left[\max_{(s,t) \in T} |a_{s,t} - f_{s,t}| \right] \leq O\left(\frac{\log^2 D}{\epsilon}\right)$$

$$\begin{aligned} \mathbb{E} \left[\max_{t \in [D]} |a_{1,t} - f_{1,t}(x)| \right] &\leq \log D \cdot O\left(\frac{\log^2 D}{\epsilon}\right), \\ &= O\left(\frac{\log^3 D}{\epsilon}\right) \end{aligned}$$

□

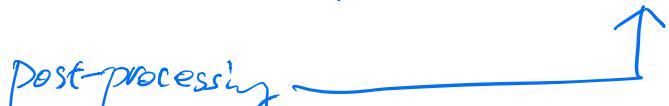
Remark:

$$① O\left(\frac{1}{\epsilon} \log^{2.5} D\right) \text{ better bound.}$$

$$\begin{aligned} ② f_{16} &= f_{14} + f_{56} \\ &= f_{18} - f_{78} \end{aligned}$$

Construct multiple estimates & Average.

(3) $f_{2,1}, f_{2,2}, \dots, f_{2,t}$ monotone.

Post-processing 

Domain Size. D

Release

$f_{s,t}(x) \quad \forall 1 \leq s \leq t \leq D.$

$$x_i \in \{1, \dots, D\}$$

$$x_i \in [0, B]$$

"Practical" solution : Discretize

$$\{0, \delta, 2\delta, \dots, B\} \quad \leftarrow \mathcal{O}\left(\frac{B}{\delta}\right).$$

↑
choosing δ is costly.

"Impossibility" Result

$$\max \text{ error} > \frac{1}{\varepsilon} \underbrace{\log^* D}_{\substack{\# \text{ take log} \\ \text{so that } D \\ \text{becomes } \leq 1}} \quad \leftarrow \log(\dots(\log(D))) \leq z.$$

Continual Release.

Stream: $x_1, x_2, \dots, x_n \in \{0,1\}$

Example: Covid Tests.

