

# Lecture 21

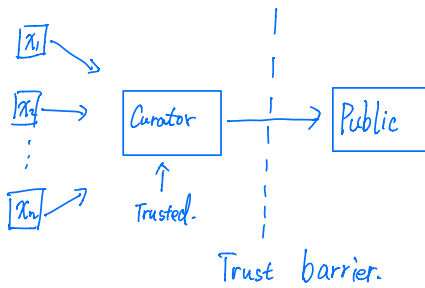
## Local Model for DP

- Frequency Estimator
- Heavy Hitters

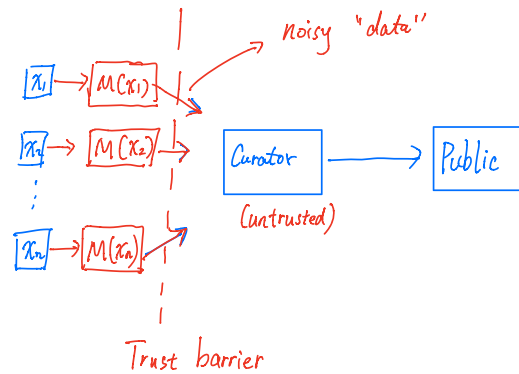
Techniques :

- Count Sketch
  - Tree Histogram
-

## Central Model



## Local Model



Local Randomizer  $M: X \mapsto Y$

is  $(\epsilon, \delta)$ -locally differentially private (LDP) if  
 $\forall x, x' \in X, E \subseteq Y$

$$\mathbb{P}[M(x) \in E] \leq e^\epsilon \mathbb{P}[M(x') \in E] + \delta.$$

## Last Lecture:

- Randomized Response
- 1-bit Mean Estimation

Error in Local Model  $\approx$  Error in Central Model  $\cdot \sqrt{n}$

$\Leftrightarrow$

↑  
quadratic.  
more samples



# Frequency Estimation in Local Model

$$x_1, \dots, x_n \in [d] \quad \{1, \dots, d\}$$

↑  
data universe.

$$f(x) = \sum_{i=1}^n \mathbb{1}[x_i = x]$$

↑  
data universe.

$$x_i \in [d] \quad x_i = \underbrace{0 \ 0 \ 0 \ \dots \ 1 \ 0 \ \dots \ 0}_{\text{length } d.}$$

$x_i$ -coordinate  
↓

Communication  
 $\Omega(d)$  bits

Runtime  $\begin{cases} \text{Server} \\ \text{User} \end{cases}$

$\Omega(d)$

↓ Local Randomizer

$$b_i = [1 \ -1 \ \dots \ -1 \ 1]$$

$$x_{ij} = 0, \quad b_{ij} = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2} \end{cases}$$

$$x_{ij} = 1, \quad b_{ij} = \begin{cases} 1 & \text{w.p. } \frac{1}{2} + \frac{\epsilon}{2} \\ -1 & \text{w.p. } \frac{1}{2} - \frac{\epsilon}{2} \end{cases}$$

$$\hat{f} = \left( \sum_i b_i \right) \cdot \frac{1}{\epsilon}$$

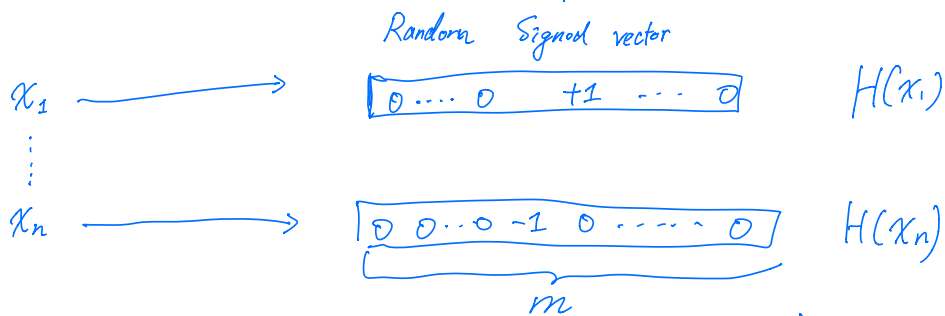
↳ error  $\approx \sqrt{n}$

Local Privacy has  $\sqrt{n}$  error

+

{ Sublinear Algorithm      has lower communication  
  Sketching                    computation  
  "Big data Algorithm"      w/  $\lesssim \sqrt{n}$  error.

# Count Sketch for dim reduction



hash function  $h: [d] \rightarrow [m] \quad (m \ll d)$   
 random sign  $s: [d] \rightarrow \{\pm 1\}$

Pairwise Independence:  $\forall x \neq x', \forall y, y'$

$$P_h [ H(x) = y \text{ and } H(x') = y' ] = \frac{1}{4m^2}$$

$\forall x \in [d]$

$$\begin{aligned} \hat{f}(x) &= \sum_i \langle H(x), H(x_i) \rangle \\ &= \sum_{i: x_i=x} \underbrace{\langle H(x), H(x_i) \rangle}_1 + \sum_{i: x_i \neq x} \langle H(x), H(x_i) \rangle \\ &= f(x) + \text{Collision Noise.} \end{aligned}$$

$$\mathbb{E}[\hat{f}(x)] = f(x) + \sum_{i: x_i \neq x} \left( \underbrace{\frac{1}{m} \left( \frac{1}{2} x 1 + \frac{1}{2} x (-1) \right)}_{h(x) = h(x_i)} + \underbrace{\left( 1 - \frac{1}{m} \right) \cdot 0}_{h(x) \neq h(x_i)} \right)$$

$$= f(x)$$

$$\downarrow$$

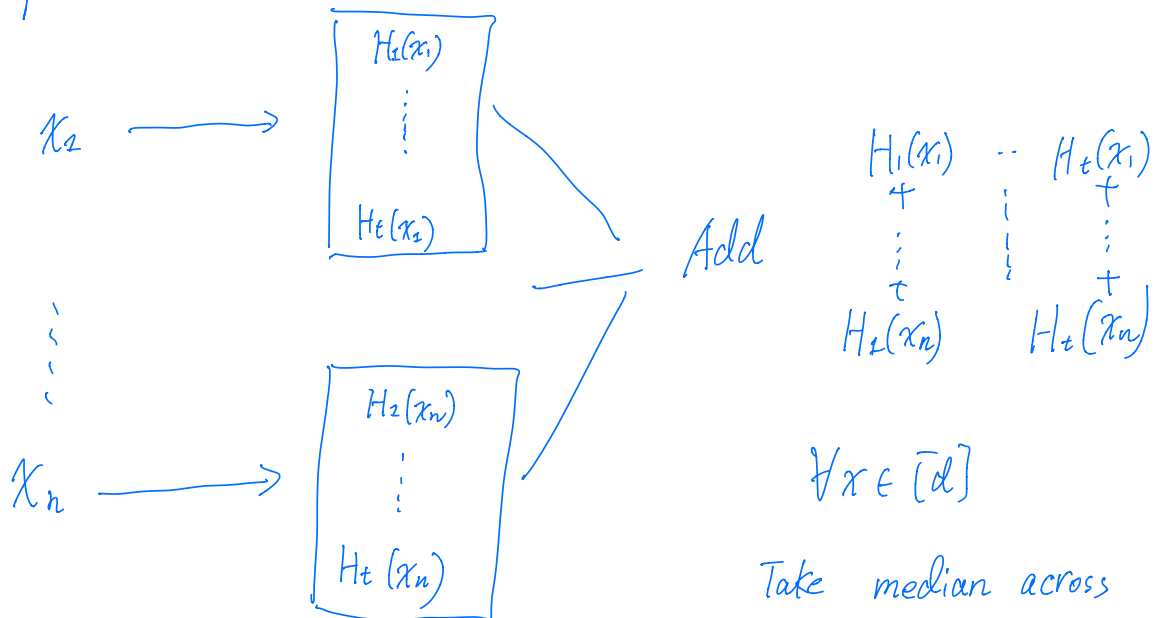
$$P[h(x) = h(x_i)] = \frac{1}{m}$$

$$P[s(x) = s(x_i)] = \frac{1}{2}$$

Set  $m \geq O(\sqrt{n})$ , w.p.  $\frac{3}{4}$

$$|\hat{f}(x) - f(x)| \leq O\left(\frac{n^{3/4}}{\sqrt{m}}\right) \rightarrow O(\sqrt{n})$$

# Amplification



Take median across the  $t$  inner products

$$\text{Set } t = \log\left(\frac{1}{\beta}\right)$$

$$\left| \hat{f}(x) - f(x) \right| \leq O\left(\frac{n^{\frac{2}{d}}}{\sqrt{m}}\right)$$

w.p.  $1 - \beta$

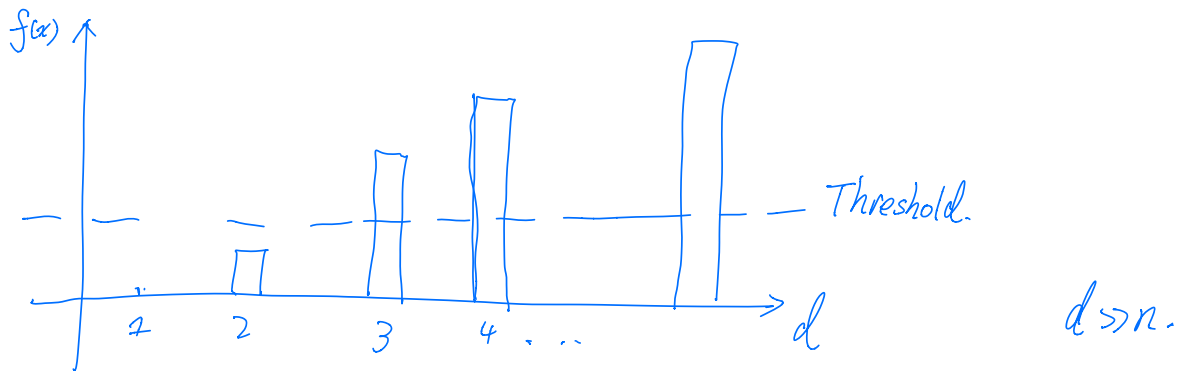
Communicate  $m = O(\sqrt{n})$  bits

Count Sketch + Random Response on each bit.



Amplification

# Heavy Hitters



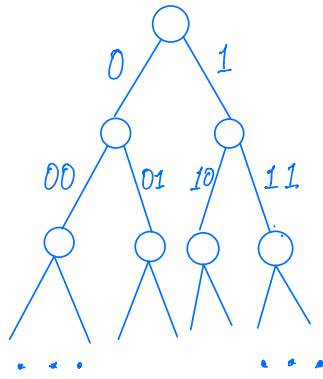
Heavy hitters  $\Rightarrow$  Frequency Estimation ?  
error  $\approx$  Threshold.

Avoid linear Runtime in  $d$  ?

TreeHist

$$x_i \in [d] \rightarrow \underbrace{(01101\dots 0)}_{\log d} \quad \mathcal{X} = \{0,1\}^{\log d}$$

Frequency oracle:  
give count for  
every sub-tree.



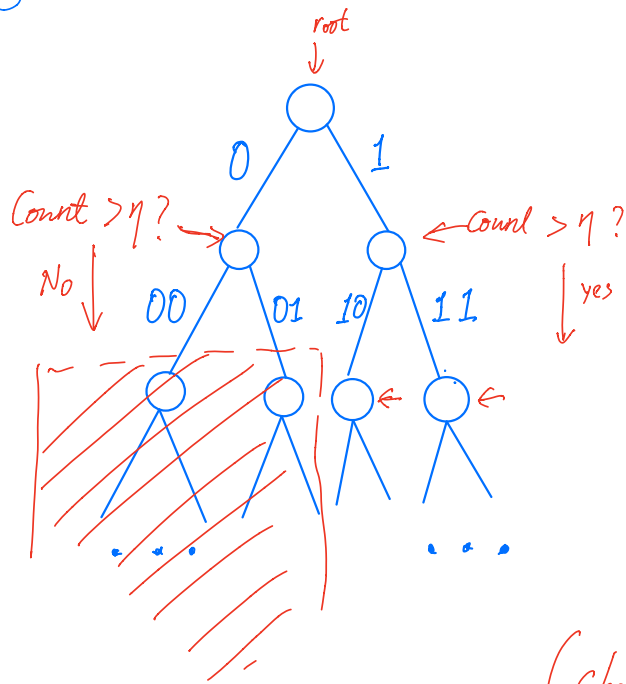
}  
 $\log d$   
Depth.

Leaf

$$\mathcal{X} = \{0,1\}^{\log(d)}$$

Problem: find  $x$   
such that  $f(x) \geq \eta$   
(e.g.  $\sqrt{n}$ )

# Pruning



Top-Down.

At each level:

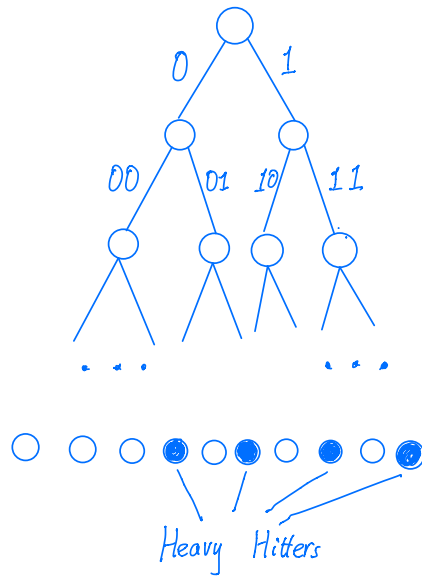
$$\# \{ \text{nodes} > \eta \} \leq \frac{n}{\eta}$$

$$\# \text{ surviving nodes at each level} \leq \frac{n}{\eta}$$

$$\# \text{ Total frequency queries} \leq O\left(\frac{n}{\eta} \log d\right).$$

(Choose  $\eta = \sqrt{n}$ )

# Heavy Hitters

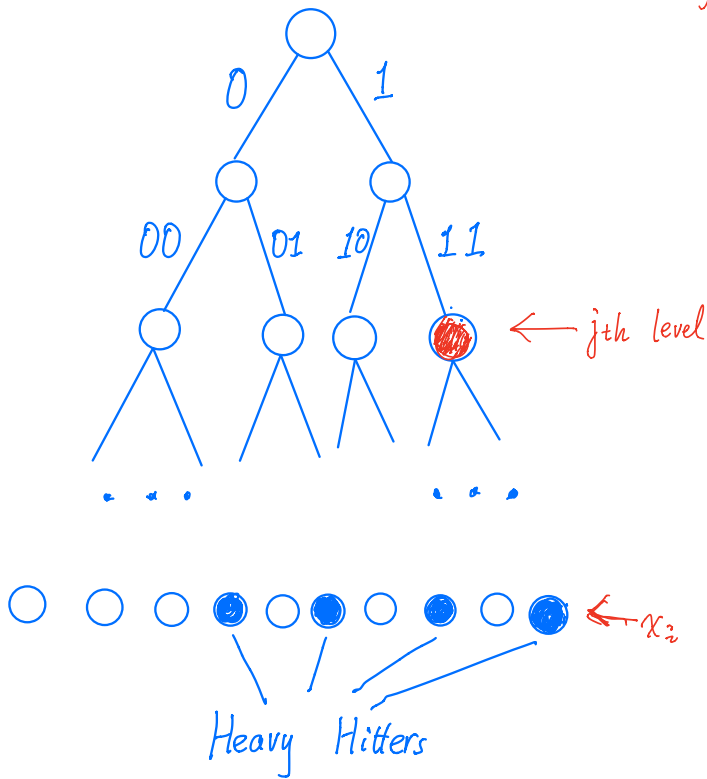


$$List = \{ (\text{heavy hitter}, \text{its frequency}) \}$$

# Frequency Estimation at each level.

$$[n] \longrightarrow [\log d]$$

Assign every  $x_i$  to a random level.



$x_i$ 's  $j$ -prefix  $\boxed{0110\dots1}$

Count Sketch  $H$

Random Signed Vector.

$\boxed{000\pm 10\dots0}$

$\epsilon$ -Randomization

$\boxed{1-1-1+1\dots-+1}$   
 $\sqrt{n}$

(f bit = 0  $\rightarrow$  unif  $\{\pm 1\}$ )

bit = 1  $\rightarrow$   $\begin{cases} +1 \text{ w.p. } \frac{1}{2} + \frac{\epsilon}{2} \\ -1 \text{ w.p. } \frac{1}{2} - \frac{\epsilon}{2} \end{cases}$

Heavy Hitters & Frequency. Pruning

Communication  $\sqrt{n}$  bits

Computation  $\sqrt{n} \log d$  for server.

• Improvement.

Hadamard Transform + Randomized  
Response.

$O(\log d)$  bits

Final error :  $O\left(\frac{\sqrt{n \log d}}{\epsilon}\right)$ .

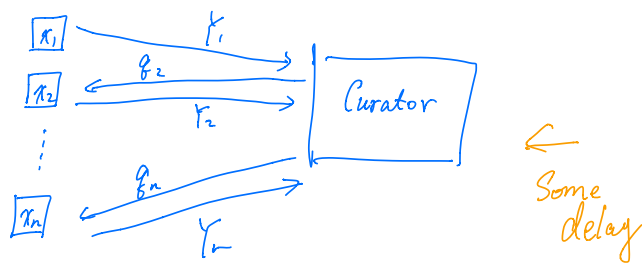
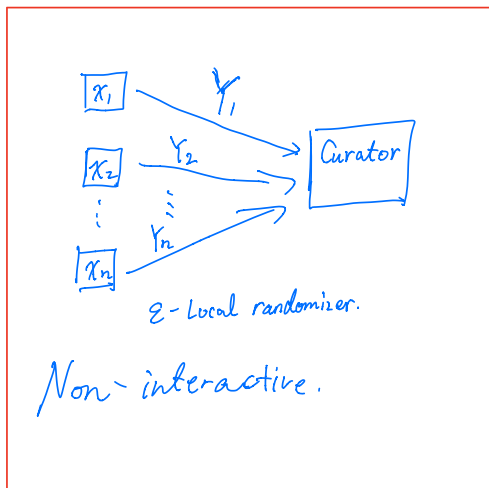
Reducing Communication Further



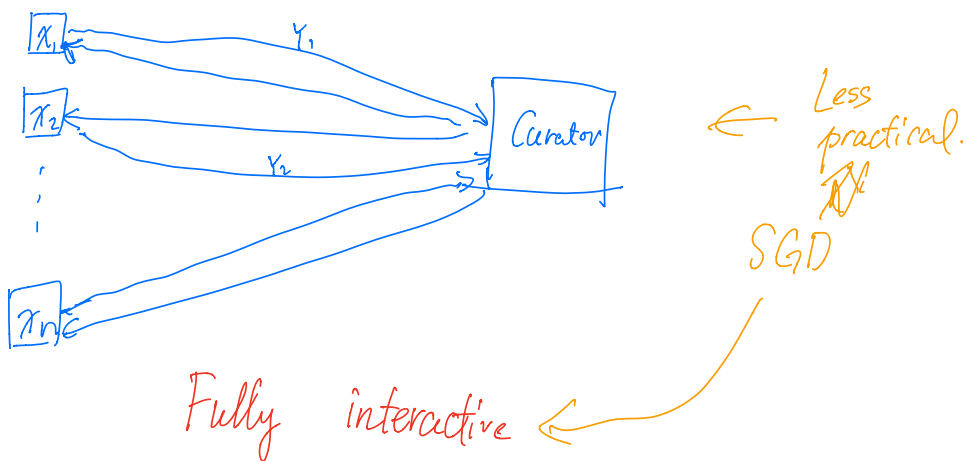


### 3 Sources of error

# Role of Interactions



Sequentially interactive  
Every user reports once.



Fully interactive  
Communication Complexity.