

Today:

- Answering Linear Queries
- Laplace & Gaussian Noise
- Special Cases
 - Histograms
 - Interval Queries
- Factorization

Toolkit so far:

→ bound global sensitivity, add noise

→ optimization, use DP-GD

Important class of problems with well-developed theory, useful all over in practice

Linear Queries \Leftrightarrow Answering Many Counts

Data $x = (x_1, \dots, x_n) \in \mathcal{U}^n$

Queries $f_j(x) = \frac{1}{n} \sum_{i=1}^n \phi_j(x_i)$

f_1, \dots, f_k

$\phi_j: \mathcal{X} \rightarrow \{0, 1\}$

"false" "true"

$F(x): \mathcal{U}^n \rightarrow [0, 1]^k$

Avg ^{Squared} error, alg A_i answers (a_1, \dots, a_k)

$$\frac{1}{k} \sum_{j=1}^k (f_j(x) - a_j)^2 = \frac{1}{k} \|F(x) - A(x)\|_2^2$$

Ask:

- ① For ϵ -DP, add $\text{Lap}(\Delta_{1/\epsilon})$ to each coordinate. What is Δ_i ?
What is avg error?

② (ϵ, δ) -DP, $\mathcal{N}(0, \Delta_2^2 \cdot \sigma_{\epsilon, \delta}^2 \cdot \Pi_d)$.

What is Δ_2 ? Avg error?

③ Apply this problem to histograms with k bins. What is Δ_1 ?

④ Let $x_i \in \mathcal{X} = [d] = \{1, 2, \dots, d\}$.

Set $\binom{d}{2}$ "threshold queries" for
 $1 \leq s \leq t \leq d$, $\phi_{s,t}(x_i) = \begin{cases} 1 & \text{if } s \leq x_i \leq t \\ 0 & \text{o.w.} \end{cases}$

What is Δ_1 ?

How can we do better?

①

$$\max_{x-x'} \|F(x) - F(x')\|_1 = \frac{k}{n}$$

$$\begin{aligned} \mathbb{E} \frac{1}{k} \|F(x) - A(x)\|_2^2 &= \frac{1}{k} \mathbb{E} \left\| \text{Lap}\left(\frac{k}{\epsilon n}\right) \right\|^2 \\ &= \frac{1}{k} \mathbb{E} \sum_{j=1}^k (\text{Lap}(\epsilon x_j))^2 \\ &= \text{Var}(\text{Lap}(\epsilon x_j)) = \frac{2k^2}{\epsilon^2 n^2} \end{aligned}$$

$$\textcircled{2} \max_{x \sim x_i} \|F(x) - F(x^*)\|_2 = \frac{\sqrt{k}}{n}$$

$$\text{cov} \Rightarrow \sigma_{\epsilon}^2 \frac{k}{n^2} \approx \frac{\log^{1/8} k}{\epsilon^2} \frac{k}{n^2} \quad \text{!!}$$

$\textcircled{3}$ ok

$$\textcircled{4} \Delta \leq d^2$$

Simple improvement :

answer subset queries, use answers to those to build up.

$$f_{i1}, f_{i2}, f_{i3}, \dots, f_{id}$$

$$\Delta_i \leq d$$

tree, take $d=8$

f_{18}

f_{14}

f_{58}

f_{12}

f_{34}

f_{56}

f_{67}

f_{11}

f_{22}

f_{33}

f_{44}

f_{55}

f_{66}

...

Claim Δ_1 for above is $\Delta_1 \leq 2 \log_2 d$.

General Question:

Given target "workload" of queries,
what's the best subset to answer?