

Lecture 9

Sparse Vector Technique

Intro

f_1, f_2, \dots, f_k

← Sens 1

Online

At time t , f_t arrives, must answer $f_t(D)$ privately

ϵ -DP. $f_t(D) + \text{Lap}(\frac{k}{\epsilon})$

(ϵ, δ) -DP. $f_t(D) + \text{Lap}(\frac{\sqrt{k}}{\epsilon})$ ← Adv. comp

poly(k) error.

Easier: which q 's are "large"?

⇒ Offline: f_1, \dots, f_k . Exp mech. return largest query.
 D dataset, objects, score = $f_i(D)$

Find c large q 's, $\frac{c \log k}{\epsilon}$ $(\frac{\sqrt{c \log k}}{\epsilon})$ approx

Online ($f_i(D) \geq T$)

- Which q 's are greater than T ?

(Goal: output first c such q 's). ← public.

Above Threshold

Privacy
 ϵ -DP

$D, D' \leftarrow$ databases

$A, A' \leftarrow$ outputs

TTTT

$a = \perp^{t-1} T$

$\Pr[A=a]$ vs $\Pr[A'=a]$ \hat{T}, v_+

Fix v_1, \dots, v_{t-1}

$g(D) = \max_{i \leq t-1} (f_i(D) + v_i)$
 deterministic

$$\Pr_{\hat{T}, v_+}[A=a] = \Pr[\hat{T} > g(D) \text{ and } f_+(D) + v_+ \geq \hat{T}]$$

$$= \Pr[g(D) \leq \hat{T} \leq f_+(D) + v_+]$$

$$= \int_{-\infty}^{\infty} \Pr[v_+ = v] \Pr[\hat{T} = \tau] \mathbb{1}_{\{\tau \in (g(D), f_+(D) + v_+)\}} d\tau dv_+$$

$$\tau = \hat{\tau} - g(D') + g(D), \quad v = \hat{v} - g(D') + g(D) - f_+(D) + f_+(D')$$

$$\Rightarrow \int \int \Pr[v_+ = v] \Pr[\hat{T} = \tau] \mathbb{1}_{\{(\hat{\tau} - g(D') + g(D)) \in (g(D), f_+(D) + \hat{v} - g(D') + g(D) - f_+(D) + f_+(D'))\}} d\hat{\tau} d\hat{v}$$

$$= \int \int \Pr[v_+ = v] \Pr[\hat{T} = \tau] \mathbb{1}_{\{\hat{\tau} \in (g(D'), \hat{v} + f_+(D'))\}} d\hat{\tau} d\hat{v}$$

$$|\tau - \hat{\tau}| \leq |g(D) - g(D')| \leq 1, \quad |v - \hat{v}| \leq 2.$$

$$\leq \int \int \exp(\epsilon/2) \Pr[v_+ = \hat{v}] \exp(\epsilon/2) \Pr[\hat{T} = \hat{\tau}] \mathbb{1}_{\{\hat{\tau} \in (g(D'), \hat{v} + f_+(D'))\}} d\hat{\tau} d\hat{v}$$

$$\exp(\epsilon) \cdot \Pr[\hat{T} \geq g(D') \text{ and } f_+(D') + v_+ \geq \hat{T}]$$

$$= e^\epsilon \Pr[A'=a].$$

Algorithm 1 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , and a threshold T . Output is a stream of responses a_1, \dots .

AboveThreshold($D, \{f_i\}, T, \epsilon$)

Let $\hat{T} = T + \text{Lap}(\frac{2}{\epsilon})$.

for Each query i do

Let $v_i = \text{Lap}(\frac{4}{\epsilon})$

if $f_i(D) + v_i \geq \hat{T}$ then

Output $a_i = \top$.

Halt.

else

Output $a_i = \perp$.

end if

end for

Accuracy

"Mistakes"

(α, β) accurate: $\forall i \in [k]$, w.p. $1 - \beta$

- If $a_i = T$, $f_i(D) \geq T - \alpha$.

- If $a_i = \perp$, $f_i(D) \leq T + \alpha$.

Thm: k queries, $\forall i < k$, $f_i(D) \leq T + \alpha$.

Then. (α, β) -accurate for $\alpha = \frac{8(\log k + \log(2/\beta))}{\epsilon}$

Proof

We compare $f_i(D) + v_i$ vs $T + (\hat{T} - T)$

Prove $\rightarrow |v_i|, |\hat{T} - T| \leq \frac{\alpha}{2}$

$a_i = T \Rightarrow f_i(D) + v_i \geq T + (\hat{T} - T)$

$f_i(D) \geq T + (\hat{T} - T) - v_i \geq T - |\alpha/2| - |\alpha/2| \geq T - \alpha$

$\Pr[|\text{Lap}(b)| \geq tb] = \exp(-t)$

$\Pr[|\hat{T} - T| \geq \frac{\alpha}{2}] = \exp(-\frac{\epsilon \alpha}{4c}) \leq \beta/2c$

$\Pr[\max_i |v_i| \geq \frac{\alpha}{2}] \leq k \cdot \exp(-\frac{\epsilon \alpha}{8c}) \leq \beta/2c$

$\log k - \frac{\epsilon \alpha}{8} \leq \log(\beta/2)$

$\hookrightarrow \log k + \log(2/\beta) \leq \frac{\epsilon \alpha}{8}$

$\alpha \geq \frac{8(\log k + \log(2/\beta))}{\epsilon}$

Sparse Vector

c outputs

Thm: (ϵ, δ) DP, $\forall \epsilon > 0, \delta \geq 0$.

Algorithm 2 Input is a private database D , an adaptively chosen stream of sensitivity 1 queries f_1, \dots , a threshold T , and a cutoff point c . Output is a stream of answers a_1, \dots .

Sparse $(D, \{f_i\}, T, c, \epsilon, \delta)$

If $\delta = 0$ Let $\sigma = \frac{2c}{\epsilon}$. Else Let $\sigma = \frac{\sqrt{32c \ln \frac{1}{\delta}}}{\epsilon}$

Let $\hat{T}_0 = T + \text{Lap}(\sigma)$

Let count = 0

for Each query i do

Let $\nu_i = \text{Lap}(2\sigma)$

if $f_i(D) + \nu_i \geq \hat{T}_{\text{count}}$ then

Output $a_i = \top$.

Let count = count + 1.

Let $\hat{T}_{\text{count}} = T + \text{Lap}(\sigma)$

else

Output $a_i = \perp$.

end if

if count $\geq c$ then

Halt.

end if

end for

$\leftarrow a_i = f_i(D) + \text{Lap}(\theta(\frac{\epsilon}{2}))$

Union bd.

Theorem 4. Suppose we are given a sequence of k queries where only c are large (i.e., the number of i such that $f_i(D) \geq T - \alpha$ is at most c). If $\delta = 0$, then Sparse is (α, β) accurate for $\alpha = \frac{8c(\log k + \log(2c/\beta))}{\epsilon}$. If $\delta > 0$, then it is (α, β) accurate for $\alpha = \frac{\sqrt{512c \log(1/\delta)}(\log k + \log(2c/\beta))}{\epsilon}$.

composition
scaling ϵ .

Numeric Sparse

- Output values of large q 's.

$$\{L, T\}^* \implies (\mathbb{R} \cup \{L\})^*$$

- If $a_i = L$, $f_i(D) \leq T + \alpha$

- If $a_i \in \mathbb{R}$, $|f_i(D) - a_i| \leq \alpha$.
 $f_i(D) \geq T - \alpha$

$$f(x) = \frac{1}{n} \sum f(x_i)$$

$$f: X^n \rightarrow [0, 1]$$

Online Private Multiplicative Weights

\hat{p} over X

Serks-1/n

True dataset p

$Q, \forall f \in Q$, estimate $f(p)$ privately

Offline:

$$\hat{p} = \text{Unif}(X)$$

1. Exp mech: choose $f \in Q$ w/ large error: $|f(\hat{p}) - f(p)|$ is large
2. Check how large $|f(\hat{p}) - f(p)|$ is. If small, return \hat{p} .
3. If large, MW update, $f(p)$ privately

Online $c = O\left(\frac{\log |X|}{\alpha^2}\right)$

$k, f_1, \dots, f_k \in Q$

Answer f_i at time i accurately

$$\hat{p} = \text{Unif}(X)$$

Run sparse vector $|f_i(p) - f_i(\hat{p})|$, $T = \Theta(\alpha)$

1. If sparse \perp : output $f_i(\hat{p})$
2. If sparse T : output $f_i(p) + \text{Lap}$
 \hookrightarrow Do an MW update using

Theorem 5. The Online Private Multiplicative Weights algorithm can, in an online manner, answer a set Q of linear queries on a database of size n to accuracy α under (ϵ, δ) -DP, given

$$n = \tilde{O}\left(\frac{\log |Q| \sqrt{\log |X| \log(1/\delta)}}{\alpha^2 \epsilon}\right) \text{ datapoints.}$$