## Sparse Vector / Above threshold:

Report Noisy Max: given $k$ queries, output one with highest value.

- Intuition: since only outputting name of one query, need less noise than if actually answering all queries.

Today: related but slightly different setting.

- Given __stream__ of queries $f_1, f_2, ...)$, each has sensitivity 1.
- Public threshold $T$
- output __first__ query above $T$, i.e., output
$$\min i: f_i(D) \geq T.$$

- Can't use exponential mechanism, since online!
- Same w/ RNM but below
- will want to generalize: output first $c$ above $T$

- Turns out to be a super useful primitive, even in non-online settings (e.g. "online" b/c iterations of larger algs).

Idea: Almost same as RNM, but also use noisy threshold.

Above Threshold:

- Let $\hat{T} = T + \text{Lap}(2/\varepsilon)$
- For each query $i$:
    - Let $\gamma_i = \text{Lap}(4/\varepsilon)$
    - if $f_i(D) + \gamma_i \geq \hat{T}$,
        - output $i$, halt

Thm: Above Threshold is $\varepsilon$-DP.

PF: Consider some $k$, $D \sim D' \in \gamma$

Fix $\gamma_1, \ldots, \gamma_{k-1}$. Take probs over $\gamma_k, \hat{T}$

$$g(D) = \max_{i < k} (f_i(D) + \gamma_i)$$

Abuse notation: $P_i[\hat{T} = t] = \text{pdf of } f_{\hat{T}} \text{ at } t.$

$$\Pr_{\hat{T}, \gamma_k}[A(D) = k] = \Pr[\hat{T} \in (g(D), f_k(D) + \gamma_k]]$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \Pr[\gamma_k = v] \cdot \Pr[\hat{T} = t] \cdot \mathbb{1}[t \in (g(D), f_k(D) + v]] \, dv \, dt$$

Change of vars: $\hat{v} = v + g(D) - g(D') + f_k(D') - f_k(D)$
$$\hat{t} = t + g(D) - g(D')$$

Note: $|\hat{v} - v| \le 2$, $|\hat{t} - t| \le 1$, since

$$g(D) - g(D') \le 1, \quad f_k(D') - f_k(D) \le 1$$

$$\text{i.e. } \hat{T} = t + g(D) - g(D')$$
$$\gamma_k = v + g(D) - g(D') + f_k(D') - f_k(D)$$

$$= \int_{-v}^{\infty} \int_{-\infty}^{\infty} P_r[\gamma_k = \hat{v}] \cdot P_r[\hat{T} = \hat{t}] \cdot$$

$$\mathbb{1}[\hat{t} \in (g(D), f_k(D) + \hat{v}]] \, dv \, dt$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P_r[\gamma_k = \hat{v}] \cdot P_r[\hat{T} = \hat{t}] \cdot$$

$$\mathbb{1}[t + g(D) - g(D') \in (g(D), v + g(D) - g(D') + f_k(D')]] \, dv \, dt$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P_r[\gamma_k = \hat{v}] \cdot P_r[\hat{T} = \hat{t}] \cdot \mathbb{1}[t \in (g(D'), f_k(D') + v]] \, dv \, dt$$

$$\underbrace{(\le_{\varepsilon_1} (\varepsilon/\varepsilon), \text{sensitivity } 2)}_{} \qquad \underbrace{(\le_{\varepsilon_1} (\varepsilon/\varepsilon), \text{sensitivity } 1)}_{}$$

$$\le \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(\varepsilon/2) P_r[\gamma_k = v] \cdot \exp(\varepsilon/2) P_r[\hat{T} = t] \cdot$$

$$\mathbb{1}[t \in (g(D'), f_k(D') + v]] \, dv \, dt$$

$$= \exp(\varepsilon) \underset{\gamma_k, \hat{T}}{P_r}[\hat{T} \in (g(D'), f_k(D') + \gamma_k]]$$

$$= \exp(\varepsilon) \cdot \underset{\hat{T}, \gamma_k}{P_r}[A(D') = k]$$

## Accuracy:

**Def:** $(\alpha, \beta)$-accurate if with prob $\geq 1-\beta$:

- Any $k$ content by alg has, $f_k(D) \geq T-\alpha$
- Any $i$ __not__ content by alg has, $f_k(D) \leq T+\alpha$

Let $\beta \in (0,1)$, and let $\alpha = \dfrac{8(\log k + \log \frac{2}{\beta})}{\varepsilon}$

**Thm:** If $f_i(D) < T-\alpha \quad \forall i < k$, then

Above Threshold is $(\alpha, \beta)$-accurate.

**Pf:** Sps $\max\limits_{i \notin C(k)} |\gamma_i| + |T - \hat{T}| \leq \alpha$

Then if we content $i$:
$$f_i(D) + \gamma_i \geq \hat{T} \Rightarrow f_i(D) \geq T - |T - \hat{T}| - |\gamma_i|$$
$$\geq T-\alpha \qquad \checkmark$$

If we don't content $i$:
$$f_i(D) + \gamma_i < \hat{T} \leq T + |T - \hat{T}|$$
$$\Rightarrow f_i(D) < T + |T - \hat{T}| + |\gamma_i| \leq T + \alpha \quad \checkmark$$

So just wts that $\max\limits_{i \notin C(k)} |\gamma_i| + |T - \hat{T}| \leq \alpha$ w.p. $\geq 1-\beta$

Laplace tail bound: If $Y \sim L_{ap}(b)$, $\Pr(|Y| \geq t \cdot b) = e^{-t}$

$$\Rightarrow Pr[\,|T-\hat{T}| \geq \tfrac{\alpha}{2}\,] = exp\left(-\tfrac{\varepsilon\alpha}{4}\right) = exp\left(-2(\log k + \log \tfrac{1}{\beta})\right)$$

$$\leq exp\left(-\log \tfrac{2}{\beta}\right) = \tfrac{\beta}{2}$$

Similarly, apply union bound to $\gamma_i$'s:

$$Pr[\max_{i\in[k]} |\gamma_i| \geq \tfrac{\alpha}{2}] \leq k \cdot exp\left(-\tfrac{\varepsilon\alpha}{8}\right) = k \cdot exp\left(-(\log k + \log \tfrac{2}{\beta})\right)$$

$$= k \cdot \tfrac{1}{k} \cdot \tfrac{\beta}{2} = \tfrac{\beta}{2} \quad \checkmark$$

Union bound $\checkmark$


Now generalize: want to output first $c$
queries above $T$.

Idea: just compose $c$ runs of Above threshold!

Sparse:
- If $\delta = 0$ let $\sigma = \tfrac{2c}{\varepsilon}$. Else let $\sigma = \dfrac{\sqrt{32 c \ln \tfrac{1}{\delta}}}{\varepsilon}$
- Let $\hat{T}_0 = T + Lap(\sigma)$
- Let count $= 0$
- For each query $i$:
    - Let $\gamma_i = Lap(2\sigma)$
    - if $f_i(D) + \gamma_i \geq \hat{T}_{count}$:
        - output i

- $count ++$
  - Let $\hat{T}_{count} = T + Lap(\sigma)$
  - if $count \geq c$ Halt.

**Thm:** Sparse is $(\varepsilon, \delta)$-DP

**Pf:** Sparse equivalent to running Above Threshold

w/ $\varepsilon' = \begin{cases} \frac{\varepsilon}{c} & \text{if } \delta = 0, \\ \dfrac{\varepsilon}{\sqrt{8c\ln\frac{1}{\delta}}} & \text{if } \delta > 0, \end{cases}$

restarting w/ fresh randomness on each output,

$\leq c$ times

If $\delta = 0$: Basic composition $\Rightarrow \leq c \cdot \varepsilon' = \varepsilon$-DP ✓

If $\delta > 0$: Advanced composition: $\sqrt{2k\ln\frac{1}{\delta}} \cdot \varepsilon' + k\varepsilon'(e^{\varepsilon'}-1)$

$\varepsilon' \Rightarrow \dfrac{\varepsilon}{\sqrt{8c\ln\frac{1}{\delta}}} \cdot \sqrt{2c\ln\frac{1}{\delta}} + c\,\dfrac{\varepsilon^2}{8c\ln\frac{1}{\delta}}$

$= \dfrac{\varepsilon}{2} + \dfrac{\varepsilon^2}{8\ln\frac{1}{\delta}} \leq \varepsilon \qquad (\varepsilon \leq 4\ln\frac{1}{\delta})$

**Accuracy:**

**Idea:** If each call to Above Threshold is
$(\alpha, \frac{\beta}{c})$-accurate, Sparse is $(\alpha, \beta)$-accurate. (union bound)

Thm: Suppose $L(T) = |\{i \leq k : f_i(D) \geq T - \alpha\}| \leq c$.

If $\delta > 0$, Sparse is $(\alpha, \beta)$-accurate for
$$\alpha = \frac{(\ln k + \ln \frac{2c}{\beta})\sqrt{512c\ln\frac{1}{\delta}}}{\varepsilon}$$

If $\delta = 0$, Sparse is $(\alpha, \beta)$-accurate for
$$\alpha = \frac{8c(\ln k + \ln \frac{2c}{\beta})}{\varepsilon}$$

Pf: Plug in accuracy for AT with
$$\varepsilon = \varepsilon', \quad \beta = \beta/c.$$
union bound.

Numeric Sparse: can also output __values__ for queries above threshold!

- Use Laplace mechanism for each, only doubles privacy loss via basic composition.

Accuracy: see book. Informally, $(\alpha, \beta)$-accurate it
w.p. $\geq (1-\beta)$, any value output within $\alpha$ of threshold,
any not output has $f_i(D) \leq T + \alpha$

Punchline: total privacy loss similar to it we
    knew which queries were above T!
    - Finding big queries is free!