# Linear Queries w/ correlated error:

- $\S(\cdot)$ class of queries $Q$ (known), each sensitivity $1$.
- Can use Laplace mechanism, composition
  - basic: noise scales like $|Q|$
  - advanced: noise scales like $\sqrt{|Q| \ln \frac{1}{\delta}}$

Basically tight: $Q$ could just be many copies of same query, need to increase noise as ask more queries!

But then trivial: return same answer each time!

In general: if $Q$ "nice", e.g., not many distinct queries, can we do better?

Goal: noise grows w/ $\log |Q|$, or better.

Offline: $Q$ known. Technique: release synthetic database that's correct on queries in $Q$. Exponential mechanism.

Online: $Q$ not known ahead of time. Can still get similar bounds. "Private multiplicative weights", sparse vector technique.

Today: offline

Setup: Linear queries.

- Generalize counting queries
- Let $X$ = domain of databases, so $P(X) = \mathcal{D}$
  (with multiplicities)

  think: every possible row of database

- Given $D \in \mathcal{D}$, $x \in X$, let $P_x$ = # copies of $x$ in $D$
- $q: X \to \{0,1\}$ predicate: "does this row correspond to having cancer?"

- Counting query: $f_q(D) = \sum_{x \in X} P_x \, q(x)$

- Normalized counting query: $f_q(D) = \frac{1}{|D|} \sum_{x \in X} P_x \, q(x)$

- <u>Linear queries</u>: $q: X \to [0,1]$
  $$f_q(D) = \sum_{x \in X} P_x \, q(x), \quad \text{or} \quad f_q(D) = \frac{1}{|D|} \sum_{x \in X} P_x \, q(x)$$

  Note: $\Delta f_q \leq 1$ unnormalized, $\Delta f_q \leq \frac{1}{|D|}$ normalized

Ex: Marginal Tables
  $(1)$ $X = \{0,1\}^d$, collection of boolean features
  (college grad, US citizen, family history of cancer...)
  Queries: what fraction of the dataset have features
  $a, b, c$?
  $S \subseteq d$, $q_s(x) = \prod_{i \in S} x_i$
  "Marginals".

$$\mathcal{Q} = \{ f_{q_S} : S \subseteq \{1, 2, \dots d\} \} \quad \text{all marginals}$$
$$|\mathcal{Q}| = 2^d$$

$$\mathcal{Q} = \{ f_{q_S} : S \subseteq [d], |S| \le k \} \quad k\text{-way marginals}$$
$$|\mathcal{Q}| = \binom{d}{k} \approx d^k$$

<u>Lots</u> of marginals, so really want to add noise
$\sim \log |\mathcal{Q}|$ rather than $\sqrt{|\mathcal{Q}|}$ or $|\mathcal{Q}|$ !


<u>Offline</u>: SmallDB $(D, \mathcal{Q}, \varepsilon, \alpha)$

  $-$ Let $R = \{ D \in \mathcal{Y} : |D| = \dfrac{\log|\mathcal{Q}|}{\alpha^2} \}$ ← small !

  $-$ Let $u : \mathcal{Y} \times R \to \mathbb{R}$ be
  $$u(D, \hat{D}) = - \max_{f \in \mathcal{Q}} |f(D) - f(\hat{D})|$$

  $-$ Use exponential mechanism $M_E(D, u, R)$ to
  sample small database from $R$


<u>Thm</u>: $\varepsilon$-DP

<u>Pf</u>: Just exponential mechanism !

<u>LtS</u>: There is a small database that's good.

<u>Lemma</u>: Let $D \in \mathcal{D}$, $Q$ collection of $\overset{(normalized)}{\text{linear queries}}$.

There exists $\hat{D}$ with $|\hat{D}| = \frac{\log |Q|}{\alpha^2}$ s.t.

$$\max_{f \in Q} |f(D) - f(\hat{D})| \leq \alpha$$

<u>Pf</u>: Let $m = \frac{\log |Q|}{\alpha^2}$.

Construct $\hat{D}$ by sampling $m$ entries from $D$ uniformly at random.

Let $Y_i$ be $i$th sample

Let $f_q \in Q$.

$$\Rightarrow f_q(\hat{D}) = \frac{1}{m} \sum_{x \in X} \hat{D}_x \, q(x) = \frac{1}{m} \sum_{i=1}^{m} q(Y_i)$$

<u>Note</u>: $0 \leq q(Y_i) \leq 1$, and

$$E[q(Y_i)] = \sum_{x \in D} \frac{1}{|D|} q(x) = f_q(D)$$

$$\Rightarrow E[f_q(\hat{D})] = \frac{1}{m} \sum_{i=1}^{m} E[q(Y_i)] = \frac{1}{m} \sum_{i=1}^{m} f_q(D) = f_q(D)$$

<u>Hoeffding bound (additive Chernoff)</u>:

Let $X_1, \ldots, X_m$ independent random vars s.t.

$0 \leq X_i \leq 1$ $\forall i$.

Thm $\Pr\left[\frac{1}{m}\sum_{i=1}^{m} X_i > E\left[\frac{1}{m}\sum_{i=1}^{m} X_i\right] + \varepsilon\right] \le \exp(-2m\varepsilon^2)$

$\underbrace{}_{E(f_q(\hat{D}))}$ $-\varepsilon\,] \overset{a}{=} \exp(-2m\,\varepsilon^2)$

$\Rightarrow \Pr\left[|f_q(\hat{D}) - f_q(D)| > \alpha\right] \le 2\,e^{-2m\alpha^2}$

union bound over all $f \in Q$:

$\Pr\left[\max_{f \in Q} |f(D) - f(\hat{D})| > \alpha\right] \le 2|Q|\,e^{-2m\alpha^2}$

$m = \dfrac{\log|Q|}{\alpha^2}$ $\Rightarrow 2|Q|\,e^{-2m\alpha^2} = 2|Q|\,e^{-2\log|Q|}$

$= \dfrac{2|Q|}{|Q|^2} < 1 \qquad (|Q| > 2)$

$\Rightarrow \exists$ good database of size $m$ $\checkmark$

So a good database exists. But we run exponential mechanism

Lemma: with prob. $\ge 1 - \beta$,

$\max_{f \in Q} |f(D) - f(\hat{D})| \le \alpha + \dfrac{2\left(\frac{\log|X|\cdot\log|Q|}{\alpha^2} + \log\frac{1}{\beta}\right)}{\varepsilon\,|D|}$

pf: use utility bound for exponential mechanism:

$$Pr\left[u(M_E(D,u,R)) \leq OPT_u(D) - \frac{2\Delta u}{\varepsilon q}\left(\log |R| + t\right)\right] \leq e^{-t}$$

$$\underset{\geq \alpha \text{ by lemma}}{\uparrow} \qquad \underset{\leq \frac{1}{|D|}}{}$$

$\Rightarrow Pr\left[\underset{f\in Q}{max} |f(D) - \hat{f}(D)| \geq -OPT_u(D) + \frac{2\Delta u}{\varepsilon}\left(\log |R| + t\right)\right] \leq e^{-t}$

2) $Pr\left[\underset{f\in Q}{max} |f(D) - \hat{f}(D)| \geq \alpha + \frac{2}{\varepsilon|D|}\left(\log\left(|X|^{\frac{(\cdot)|Q|}{\alpha^2}}\right) + \ln \frac{1}{\beta}\right)\right] \leq \beta$

3) $Pr\left[\underset{f\in Q}{max} |f(D) - f(\hat{D})| \geq \alpha + \frac{2}{\varepsilon|D|}\left(\frac{(\cdot)|Q|\log|X|}{\alpha^2} + \ln \frac{1}{\beta}\right)\right] \leq \beta$ ✓


Thm: For any database $D$ with

$$|D| \geq \frac{16 \log |X| \log |Q| + 4\alpha^2 \log \frac{1}{\beta}}{\varepsilon \alpha^3}$$

w.p. $\geq 1-\beta$, $\underset{f\in Q}{max} |f(D) - f(\hat{D})| \leq \alpha$

pf: use previous lemma w/ $\frac{\alpha}{2}$: w.p. $\geq 1-\beta$

$$\underset{f\in Q}{max} |f(D) - f(\hat{D})| \leq \frac{\alpha}{2} + \frac{2\left(\frac{4 (\cdot)|X| \log |Q|}{\alpha^2} + \log \frac{1}{\beta}\right)}{\varepsilon |D|}$$

set to be $\leq \alpha$, solve for $|D|$:

$$\frac{\alpha}{2} + \frac{2\left(\frac{4 (\cdot)|X| \log |Q|}{\alpha^2} + \log \frac{1}{\beta}\right)}{\varepsilon |D|} \leq \alpha$$

$$\Longleftrightarrow \frac{\alpha}{2} |D| + \frac{1}{\varepsilon}\left( \frac{8 \log |\mathcal{X}| \log |Q|}{\alpha^2} \, 2\log \frac{1}{\beta} \right) \leq \alpha |D|$$

$$\Longleftrightarrow |D| \geq \frac{1}{\varepsilon}\left( \frac{16 \log |\mathcal{X}| \log |Q| + 4\alpha^2 \log \frac{1}{\beta}}{\alpha^3} \right)$$

interpretation: Think of a small constant, e.s., $\frac{1}{10}$

Can answer all queries w/ database of size
$$\sim \frac{1}{\varepsilon} \log |\mathcal{X}| \log |Q| \; !$$

e.g. all marginals w/ $\sim \frac{1}{\varepsilon} d^2 \; !$

More refined bounds:

VC-dimension of $Q$ replaces $\log |Q|$ for counting queries

fat-shattering dim of $Q$ for linear queries.