# (Linear) Query Release & Synthetic Data

- Recap    query release problem
- Why    " linear " ?
- Factorization   Framework.
    $\hookrightarrow$   " Matrix Mechanism "

# Linear Query Release

Dataset $\quad x = (x_1, \ldots, x_n) \in \mathcal{X}^n$

$\hookleftarrow$ "data universe"

Statistics $\quad f_1, \ldots, f_k$

$$f_i(x) = \frac{1}{n} \sum_{j=1}^{n} c_i(x_j) \quad, \quad c_i : \mathcal{X} \mapsto \{0,1\}$$

"predicate"

Asian $\wedge$ "age $\geq 30$" $\to 1$

Goal: output $\vec{a} = (a_1, \ldots, a_k)$

$$\left( \frac{1}{k} \sum_{i=1}^{k} \left( f_i(x) - a_i \right)^2 \right)^{\frac{1}{2}} \leq \alpha \qquad \left( \begin{array}{l} \text{Extends to} \\ c_i : \mathcal{X} \mapsto \mathbb{R} \end{array} \right)$$

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{"}\ell_2 \text{ error"}}$

$$\vec{F}(x) = (f_1(x), \ldots, f_k(x))$$

$$\frac{1}{\sqrt{k}} \ \| \vec{F}(x) - a \|_2 \quad \leq \alpha.$$

Mechanisms: Laplace
Gaussian
Binary Tree Mechanism

# Gaussian Mechanism

$$M(x) = F(x) + Z \quad , \quad Z \sim N(0, \sigma^2 I_{k \times k}).$$

$$\sigma^2 = C_{\varepsilon, \delta} \cdot \Delta_2^2$$

$$\approx \frac{\log(1/\delta)}{\varepsilon^2} \qquad \hookleftarrow \ell_2\text{-sensitivity}$$

$$\Delta_2 = \max_{x \sim x'} \| F(x) - F(x') \|_2$$

$$= \max_{u, u' \in \mathcal{X}} \frac{1}{n} \| F(u) - F(u') \|_2$$

$$\leq \max_{u \in \mathcal{X}} \frac{2}{n} \| F(u) \|_2$$

$$\leq 2 \frac{\sqrt{k}}{n}$$

$\ell_2$ error
$$\mathbb{E}\left[ \frac{\| Z \|_2}{\sqrt{k}} \right] \leq C_{\varepsilon, \delta} \cdot \sigma$$

$$= \frac{2 \cdot C_{\varepsilon, \delta} \ k^{\frac{1}{2}}}{n}$$

# Why "Linear"?

dataset $\quad x = (x_1, \ldots, x_n) \in \mathcal{X}^n$

data universe $\quad \mathcal{X} = \{1, \ldots, m\} \quad , \quad |\mathcal{X}| = m$

Histogram $\quad h_x \in \mathbb{R}^m$

$$\forall u \in \mathcal{X} : (h_x)_u = \frac{1}{n} \left| \{ j : x_j = u \} \right|$$

univ: $\mathcal{X} = \{1, 2, 3\}$

dataset: $x = (1, 2, 3, 3.1) \quad \longrightarrow \quad h_x = \left( \frac{2}{5}, \frac{1}{5}, \frac{2}{5} \right)$

$h_{x'} = \left( \frac{1}{5}, \frac{2}{5}, \frac{2}{5} \right)$

① $\forall$ dataset $x$, $\quad \| h_x \|_1 = 1$

② $\forall$ neighbors $x \& x'$, $\quad \| h_x - h_{x'} \|_1 \leq \frac{2}{n}$

③ Queries $f_1, \ldots, f_k$

$\quad \varphi_1, \ldots, \varphi_k$

$$F = \quad k \begin{pmatrix} \varphi_1(u_1) & \cdots & & \varphi_1(u_m) \\ \vdots & \ddots & \ddots & \vdots \\ \varphi_k(u_1) & \cdots & \cdots & \varphi_k(u_m) \end{pmatrix} \quad m$$

answer vector $\quad \vec{F}(x) = F h_x$

"Only for thought experiment" $\qquad m \neq 2^d$

# Revisiting Gaussian Mechanism.

$$M(x) = F h_x + Z \qquad , \qquad Z \sim N(0, \beta^2 I_{k \times k})$$

$$\beta^2 = C_{\varepsilon,\delta}^2 \cdot \underline{\Delta_2^2}$$

$$\Delta_2 = \max_{x \sim x'} \| F \underbrace{(h_x - h_{x'})}_{v} \|_2$$

$$\leq \max_{\substack{v \in \mathbb{R}^m \\ \|v\|_2 = \frac{2}{n}}} \| F v \|_2$$

$$= \frac{2}{n} \underbrace{\max_{\substack{v \in \mathbb{R}^m \\ \|v\|_1 = 1}} \| F v \|_2}_{} \longrightarrow \text{" largest } \ell_2\text{-norm of}$$
$$\text{a column in } F\text{"}$$

$$= \frac{2}{F} \quad \| F \|_{1 \to 2}$$

$$\mathbb{E}\left[ \frac{\| Z \|_2}{\sqrt{k}} \right] = 2 \quad \underline{\frac{C_{\varepsilon,\delta} \cdot \boxed{\| F \|_{1 \to 2}}}{n}} \qquad \text{Characterize Sensitivity.}$$

$$\text{worst-case } \sqrt{k}$$

Can we do better when

F has some structure ?

# General Factorization Framework

Histogram $h_x \in \mathbb{R}^m$ $\begin{cases} \text{dataset of size } n \\ \text{linear queries} \quad F \in \mathbb{R}^{k \times m} \end{cases}$

Want to release $F h_x$

① Approximate $\tilde{F} \approx F$

② Factorize $\tilde{F} = R \quad M$

    "Reconstruction" ↑    ↳ "measurement"

$$\hat{a} = R \left( M h_x + \underset{\uparrow}{z} \right)$$

           ↳ noise

$$= R M h_x + R z$$
$$= \tilde{F} h_x + R z$$

③ Post-processing to $\tilde{a}$
   to satisfy some "consistency" properties

*Focus of today*

# Factorization

Given linear queries $F \in \mathbb{R}^{k \times m}$

(Trivial Example) $\qquad f_1 = \cdots = f_k$

Gaussian Mechanism: error $\approx \dfrac{k^{\frac{1}{2}}}{n} \cdot C_{\varepsilon, \delta}$

$$a_1 = f_1(x) + z$$

$$\vec{a} = (a_1, \ldots, a_1) \qquad \text{error} \approx \frac{1}{n} \; C_{\varepsilon, \delta}$$

$$F = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ & & \cdots & & \\ 1 & 0 & 1 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 1 & 0 & 1 \end{pmatrix}$$

$$R \in \mathbb{R}^{k \times \ell} \qquad\qquad M \in \mathbb{R}^{\ell \times m}$$

More generally:

Gaussian Mech $\qquad$ v.s.

$F h_x + \boxed{z}$   i.i.d noise.

$\propto \|M\|_{1 \to 2}$

$$R(M h_x + z)$$
$$= F h_x + \boxed{R z}$$

proportional to
$\|F\|_{1 \to 2}$ $\quad \begin{pmatrix} \text{could be} \\ \sqrt{k} \end{pmatrix}$

Correlated noise.

# Factorization

For $R, M$ s.t. $F = RM$

$$M_{R,M}(x) = R\left(M h_x + Z\right) \quad \rightarrow \text{error}$$

$$= F h_x + \boxed{RZ}$$

$\leftarrow$ Correlated noise.

$$Z \sim N\left(0, \, \beta^2 I_{\ell \times \ell}\right)$$

$$\beta^2 = C_{\varepsilon, \delta}^2 \, \|M\|_{1 \rightarrow 2}^2$$

# Factorization

For fixed $RM = F$, analyze error $\mathbb{E}\left[\dfrac{\|RZ\|_2}{\sqrt{k}}\right]$

$$\underbrace{\begin{pmatrix} r_1 \cdot z \\ \vdots \\ r_k \cdot z \end{pmatrix}}_{RZ} = \underbrace{\begin{pmatrix} \rule{1.5em}{0.4pt}\, r_1 \,\rule{1.5em}{0.4pt} \\ \vdots \\ \rule{1.5em}{0.4pt}\, r_k \,\rule{1.5em}{0.4pt} \end{pmatrix}}_{R} \underbrace{\begin{pmatrix} z_1 \\ \vdots \\ z_\ell \end{pmatrix}}_{Z}, \qquad Z \sim N(0, \delta^2 I_{\ell \times \ell})$$

---

**Fact.** $\quad r_i \cdot z \sim N\left(0, \delta^2 \|r_i\|_2^2\right)$

$\qquad \mathbb{E}\left[(r_i \cdot z)^2\right] = \delta^2 \|r_i\|_2^2 .$

---

$\mathbb{E}\left[\|RZ\|_2\right] \underset{\text{Jensen}}{\leq} \left(\mathbb{E}\left[\|RZ\|_2^2\right]\right)^{\frac{1}{2}}$

$\qquad = \left(\mathbb{E}\left[\sum_{i=1}^{k} (r_i \cdot z)^2\right]\right)^{\frac{1}{2}}$

$\qquad = \left(\sum_{i=1}^{k} \mathbb{E}\left[(r_i \cdot z)^2\right]\right)^{\frac{1}{2}}$

$$\begin{pmatrix} \rule{1.5em}{0.4pt}\, r_1 \,\rule{1.5em}{0.4pt} \\ \vdots \\ \rule{1.5em}{0.4pt}\, r_k \,\rule{1.5em}{0.4pt} \end{pmatrix}$$
$R$

$\qquad = \left(\sum_{i=1}^{k} \delta^2 \|r_i\|_2^2\right)^{\frac{1}{2}}$

$\qquad = \delta \left(\underbrace{\sum_{i=1}^{k} \|r_i\|_2^2}\right)^{\frac{1}{2}}$

$\qquad = \underset{\substack{\uparrow \\ \text{Scales with} \\ \|M\|_{1\to 2}.}}{\delta} \quad \underset{\substack{\uparrow \\ \longrightarrow \text{Frobenius norm.}}}{\|R\|_F}$

Putting together : $\qquad$ Expected error

$$\frac{1}{\sqrt{k}} \mathbb{E}\left[\|RZ\|_2\right] \leq O\left(\frac{C_{\varepsilon,\delta} \quad \|R\|_F \quad \|M\|_{1\to 2}}{\sqrt{k} \qquad n}\right)$$

Example / Exercise:

Binary Tree Mechanism.

# Factorization Framework.

Error
$$O\left(\frac{C_{\varepsilon,\delta}}{n} \cdot \frac{\|R\|_F \cdot \|M\|_{1\to 2}}{\sqrt{k}}\right)$$

Factorization norm of $F$

$$\gamma(F) = \min\left\{\frac{\|R\|_F \cdot \|M\|_{1\to 2}}{\sqrt{k}} : RM = F\right\}$$

**Theorem.** For every $F \in \mathbb{R}^{k\times m}$, there is $(\varepsilon, \delta)$- DP mechanism with $\ell_2$ - error $\leq O\left(\frac{C_{\varepsilon,\delta}}{n} \cdot \gamma(F)\right)$.