

(Linear) Query Release

- Recap : Factorization Mechanism
- Projection Mechanism
- Preview : Synthetic Data
- Online Learning

Update : lecture notes / Slides

Linear Query Release

Dataset $\mathcal{X} = (x_1, \dots, x_n) \in \mathcal{X}^n$
"data universe"

Statistics f_1, \dots, f_k

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

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

$$\left(\frac{1}{k} \sum_{i=1}^k (f_i(x) - a_i)^2 \right)^{\frac{1}{2}} \leq \alpha$$

" l_2 error"

$$\vec{F}(x) = (f_1(x), \dots, f_k(x))$$
$$\frac{1}{\sqrt{k}} \|\vec{F}(x) - \vec{a}\|_2 \leq \alpha.$$

Why "Linear"?

dataset $x = (x_1, \dots, x_n) \in \mathcal{X}^n$
data universe $\mathcal{X} = \{1, \dots, m\}$, $|\mathcal{X}| = m$

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

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

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

Queries f_1, \dots, f_k
 $\varphi_1, \dots, \varphi_k$

$$F = \begin{matrix} & & & m \\ & & & \\ & & & \\ & & & \\ & & & \end{matrix} \begin{pmatrix} \varphi_1(u_1) & \dots & \varphi_1(u_m) \\ \vdots & & \vdots \\ \varphi_k(u_1) & \dots & \varphi_k(u_m) \end{pmatrix}$$

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

General Factorization Framework

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

{ dataset of size n
linear queries $F \in \mathbb{R}^{k \times m}$

Want to release $F h_x$

① Approximate $\tilde{F} \approx F$

② Factorize $\tilde{F} = R \quad M$
"Reconstruction" "measurement"

$$\begin{aligned}\hat{a} &= R (M h_x + Z) \\ &= R M h_x + R Z \\ &= \tilde{F} h_x + R Z\end{aligned}$$

noise

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

Factorization

$$\text{For } R, M \quad \text{s.t.} \quad F = \overset{\text{Reconstruction}}{\downarrow} R \overset{\text{Measurement}}{\downarrow} M$$

$$\begin{aligned} M_{R,M}(x) &= R \left(\underbrace{M h_x}_{\text{Reconstruction}} + Z \right) \\ &= F h_x + \underbrace{RZ}_{\text{Correlated noise}} \end{aligned}$$

$$Z \sim N(0, \sigma^2 I_{l \times l})$$

$$\sigma^2 = c_{e,s}^2 \|M\|_{1 \rightarrow 2}^2$$

Factorization Framework.

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

Factorization norm of F

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

Theorem. For every $F \in \mathbb{R}^{k \times m}$, there is (ϵ, δ) -DP mechanism

$$\text{with } \ell_2\text{-error} \leq O\left(\frac{C_{\epsilon, \delta}}{n} \cdot \chi(F)\right).$$

Standard
Gaussian

$$O\left(\frac{C_{\epsilon, \delta}}{n} \cdot \sqrt{k}\right)$$

→ Factorization \approx preprocessing.
Gaussian

→ post-processing: Projection mechanism

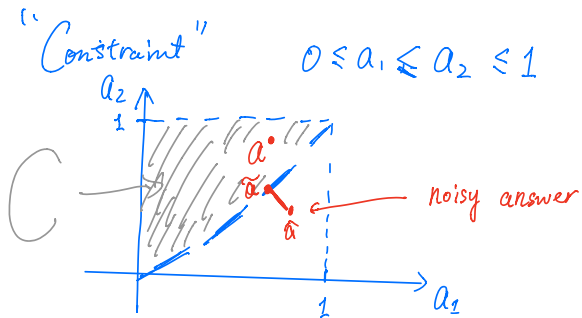
Consistency

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

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

$$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = F \cdot h_x$$

$$\text{Gaussian: } \begin{pmatrix} \hat{a}_1 \\ \hat{a}_2 \end{pmatrix} = F h_x + z$$



Projection:

$$\tilde{a} = \Pi_C(\hat{a}) = \arg \min_{a' \in C} \|a' - \hat{a}\|_2$$

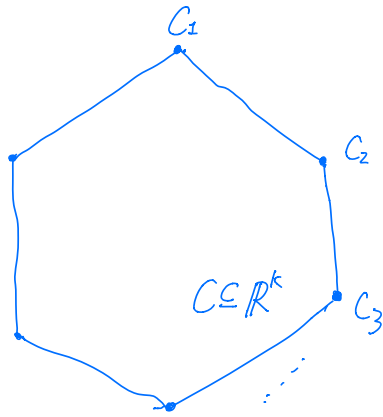
- ① \tilde{a} is consistent
- ② $\|\tilde{a} - a\|_2 \leq \|\hat{a} - a\|_2$
- ③ \tilde{a} is (ϵ, δ) -DP.

Projection may help accuracy!

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

$$a_1 + a_2 = a_3$$

Consistency



"Convex hull"

Queries matrix

$$F = \begin{matrix} & & & m \\ & & & \\ k & \left(\begin{array}{c|c|c|c} & & & \\ \hline c_1 & c_2 & \dots & c_m \\ \hline & & & \end{array} \right) & & \end{matrix}$$

$$F = \begin{pmatrix} \varphi_1(u_1) & \dots & \varphi_1(u_m) \\ \vdots & & \vdots \\ \varphi_k(u_1) & \dots & \varphi_k(u_m) \end{pmatrix}$$

$$C = \{ a \in \mathbb{R}^k : \exists h \in \mathbb{R}_+^m, \|h\|_1 = 1, a = Fh \}$$

The Projection Mechanism

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

Gaussian Mechanism :

$$a = Fh$$

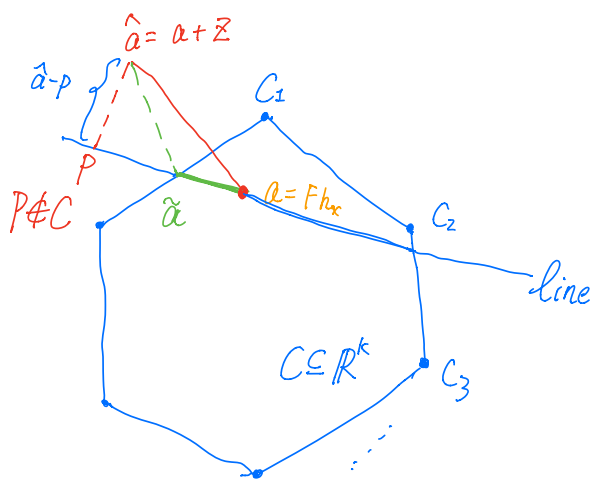
$$\hat{a} = a + Z, \quad Z \sim \mathcal{N}(0, \sigma^2 I_{k \times k})$$

$$\sigma^2 = \frac{C_{\epsilon, \delta}^2 k}{n^2}$$

Projection

$$\text{return } \tilde{a} = \arg \min_{a' \in C} \|a - a'\|_2$$

Analysis



" l_2 error"

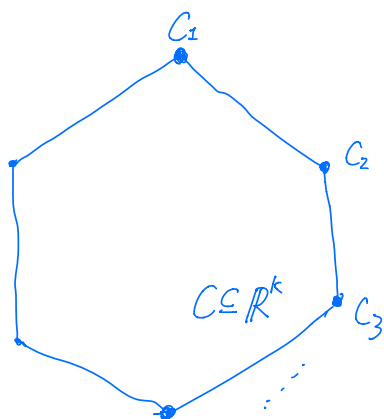
$$\begin{aligned}
 & \| \tilde{a} - a \|_2^2 \\
 &= \langle \tilde{a} - a, \tilde{a} - a \rangle \\
 &\leq \langle \tilde{a} - a, p - a \rangle \\
 &= \langle \tilde{a} - a, p - a + \underbrace{\hat{a} - p}_{\perp \tilde{a} - a} \rangle \\
 &= \langle \tilde{a} - a, \hat{a} - a \rangle \\
 &= \langle \tilde{a} - a, z \rangle \\
 &= \langle \tilde{a}, z \rangle - \langle a, z \rangle \\
 &\leq |\langle \tilde{a}, z \rangle| + |\langle a, z \rangle| \\
 &\leq \underbrace{2 \max_{v \in C} |\langle v, z \rangle|}_{\text{Data independent.}}
 \end{aligned}$$

Analysis .

$$\mathbb{E} \left[\frac{\|\tilde{a} - a\|_2}{\sqrt{k}} \right] \stackrel{\text{Jensen Ineq}}{\leq} \mathbb{E} \left[\frac{\|\tilde{a} - a\|_2^2}{k} \right]^{\frac{1}{2}}$$



$$\text{(Plug in)} = \left(\frac{2 \mathbb{E}_Z \left[\max_{v \in C} |\langle v, Z \rangle| \right]}{k} \right)^{\frac{1}{2}}$$



Fact. (Suffices to think about vertices)

$$\max_{v \in C} |\langle v, Z \rangle| = \max_{j \in [m]} |\langle C_j, Z \rangle|$$

Columns of F

$$\langle C_j, Z \rangle \sim \mathcal{N}(0, \|C_j\|_2^2 \sigma^2)$$

Fact. If W_1, \dots, W_m are Gaussian with variance $\leq (\sigma')^2$, & zero-mean then $\mathbb{E} \left[\max \{|W_1|, \dots, |W_m|\} \right] \leq \sigma' \cdot \sqrt{\log m}$

$$\begin{aligned} \Rightarrow \mathbb{E}_Z \left[\max_{v \in C} |\langle v, Z \rangle| \right] &= \mathbb{E}_Z \left[\max_{C_j, j \in [m]} |\langle C_j, Z \rangle| \right] \\ &= \sqrt{k} \cdot \sigma \cdot \sqrt{\log m}. \end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[\frac{\|\tilde{a} - a\|_2}{\sqrt{k}} \right] &\stackrel{\text{Jensen Ineq}}{\leq} \mathbb{E} \left[\frac{\|\tilde{a} - a\|_2^2}{k} \right]^{\frac{1}{2}} \\
&= \left(\frac{2 \mathbb{E} \left[\max_{v \in \mathcal{C}} \langle v, z \rangle \right]}{k} \right)^{\frac{1}{2}} \\
&\leq \mathcal{O} \left(\left(\frac{C_{\epsilon, \delta} \cdot \sqrt{\log m}}{n} \right)^{\frac{1}{2}} \right).
\end{aligned}$$

Projection Mechanism
Bound.

Just
Gaussian Mechanism

$$\min \left\{ \mathcal{O} \left(\left(\frac{C_{\epsilon, \delta} \cdot \sqrt{\log m}}{n} \right)^{\frac{1}{2}} \right), \mathcal{O} \left(\frac{C_{\epsilon, \delta} \cdot \sqrt{k}}{n \epsilon} \right) \right\}$$

Dependence on n

$$\frac{1}{\sqrt{n}} \quad \text{v.s.} \quad \frac{1}{n}$$

Dependence on k

$$\text{No dep.} \quad \text{v.s.} \quad \sqrt{k}$$

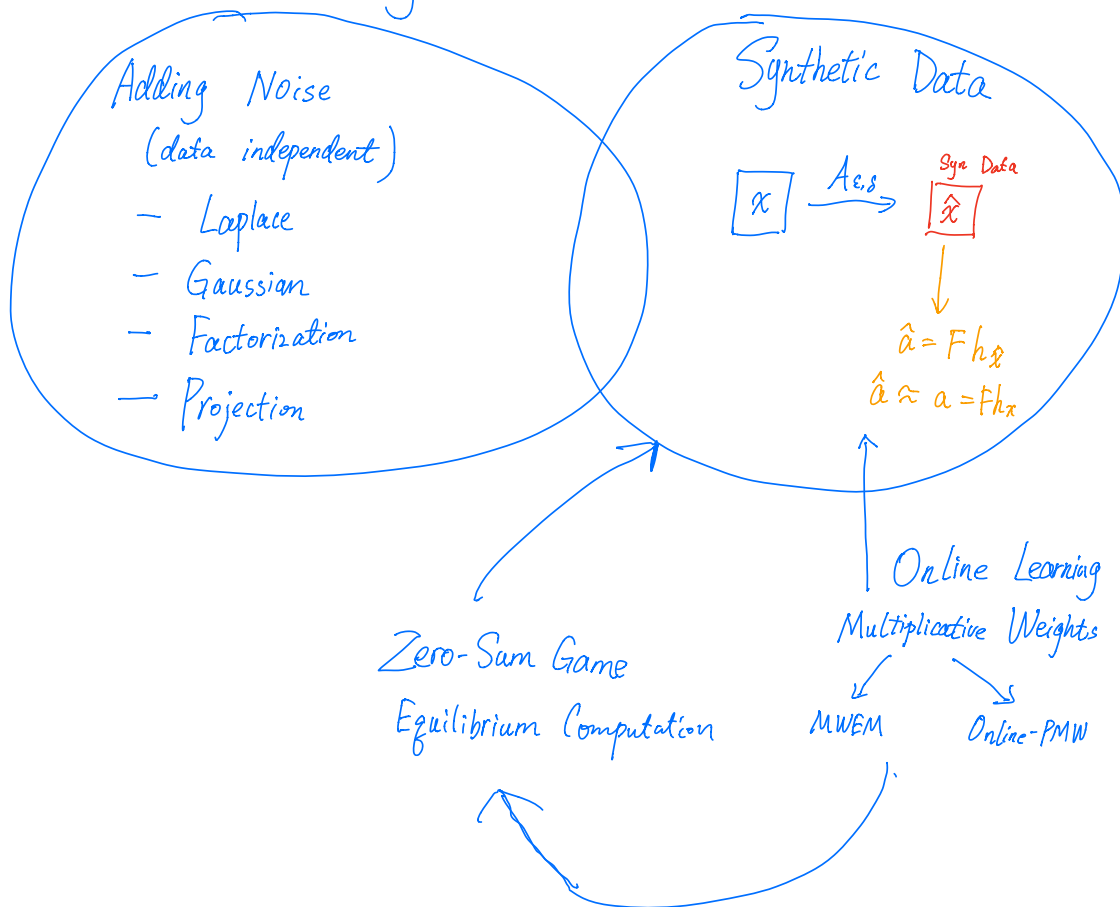
Dependence on m

$$\sqrt{\log m} \quad \text{v.s.} \quad \text{No dep.}$$

Example: $\mathcal{X} = \{0, 1\}^d$, $m = 2^d$, $k = 2^d$

$$\left(\frac{d^{\frac{1}{2}}}{n} \right)^{\frac{1}{2}} \quad \text{v.s.} \quad \frac{2^{d/2}}{n}$$

Query Release



Online Learning

"Sequential decision making"

- Setting:
- a set of actions $\{1, \dots, k\} = [k]$.
 - "Game" between decision-maker D & Adversary A .

"may know
 D 's alg
not D 's randomness"

For $t=1, \dots, T$:

D chooses distribution $p^t \in \Delta(A)$

A chooses cost vector $c^t \in [0, 1]^k$

$a^t \sim p^t$ sampled action

D pays cost $C_{a^t}^t$ and observes c^t .

Focus on:

$$\mathbb{E}_{a \sim p^t} [C_a^t] = \langle p^t, c^t \rangle$$

$$\text{Total cost} = \sum_{t=1}^T C_{a^t}^t$$

How to measure D ?

— Compare w/ the sequence $a^1 \dots a^T$

Hopeless

— Compare w/ the best action in hindsight

$$\text{Regret} = \frac{1}{T} \sum_{t=1}^T C_{a^t}^t - \min_{a \in [k]} \frac{1}{T} \sum_{t=1}^T C_a^t$$

Online Learning

Give an algorithm *Multiplicative Weights*.

$$\mathbb{E}[\text{Regret}] \leq O\left(\sqrt{\frac{\ln(k)}{T}}\right).$$