

## Lecture 24

Adaptive Data Analysis. (ADA)

- Reusable Holdout
- Sparse Vector Mechanism (SVM)

↳ Application: Synthetic Data for ADA.

# Logistics

Today is the last day of "me" lecturing.

Weds : Additional Office Hour for Projects.

Fri : Continue office hour

Next Week : Project Presentation

Schedule also on Canvas

= May 3rd

Ryan Steed

Shengyuan Hu

Justin Whitehouse

Tianshi Li

= May 5th

Charlie Hou

Zhili Feng

Shuaiqi Wang

Format :

Expect Slides

20 minutes including Q&A.

Final Write-up due later.

## Model for Adaptive Data Analysis

### Statistical / Linear Queries

$$\phi: \mathcal{X} \rightarrow [0,1] \quad \text{"predicate"}$$

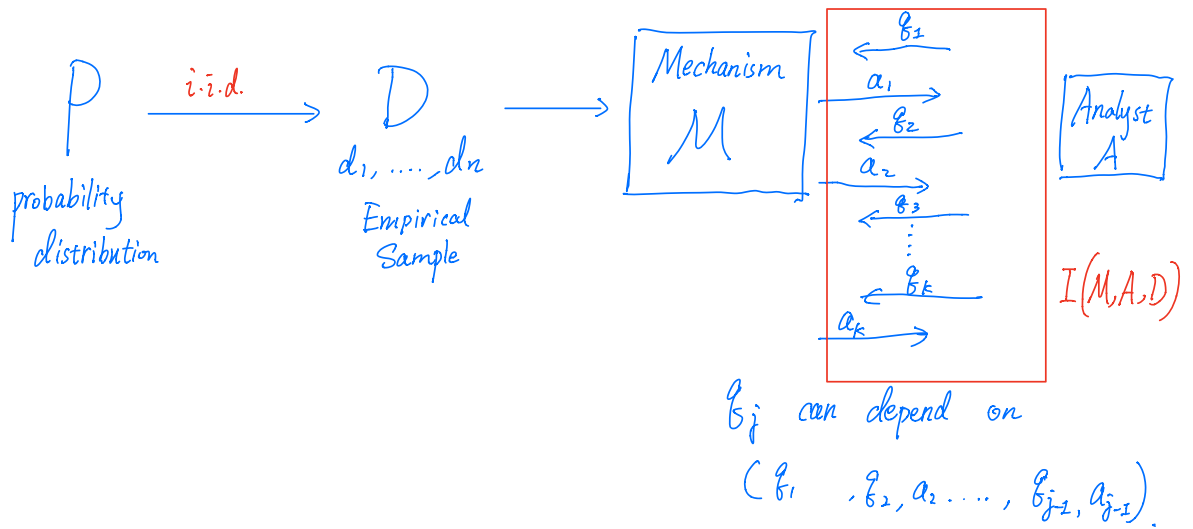
$$q_\phi(P) = \mathbb{E}_{x \sim P} [\phi(x)] \quad \text{"Population value"}$$

↑  
population

$$q_\phi(D) = \hat{\mathbb{E}}_{d_i \sim D} [\phi(d_i)] = \frac{1}{n} \sum_{i=1}^n [\phi(d_i)] \quad \text{"Empirical Average"}$$

$$D = (d_1, \dots, d_n) \in \mathcal{X}^n$$

# Interaction of Adaptive Data Analysis.



Transfer Theorem.  $(\epsilon, \delta)$ -version [JLNRSS20]

Suppose  $I(M, A, D)$  is  $(\alpha, \beta)$ -sample accurate  $\leftarrow$   
 &  $(\epsilon, \delta)$ -DP.  $\leftarrow$

Then for every  $c, d > 0$ ,  $I(M, A, D)$  is  $(\alpha', \beta')$   $\leftarrow$   
 distributionally accurate, for

$$\alpha' = \alpha + \underbrace{(e^\epsilon - 1)}_{\approx \epsilon} + \underbrace{c}_{\alpha} + \underbrace{2d}_{\epsilon}, \quad \beta' = \frac{\beta}{c} + \frac{\delta}{d}$$

$$\alpha' = O(\alpha + \epsilon) \quad \beta' = \left( \frac{\beta}{\alpha} + \frac{\delta}{\epsilon} \right)$$

Answering  $k$  adaptive queries.

Gaussian Mechanism v.s. Sample Splitting

$$\alpha'' \approx \sqrt{\frac{k^2}{n}}$$

$\uparrow$   
 (can answer  $k \approx O(n^2)$  queries)

$$\alpha' \leq \sqrt{\frac{k}{n}}$$

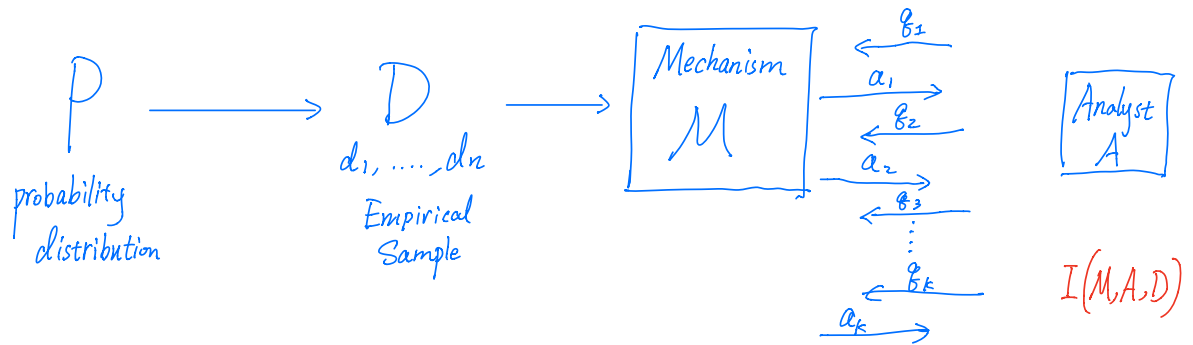
$\uparrow$   
 $k \leq O(n)$

Can we do better?

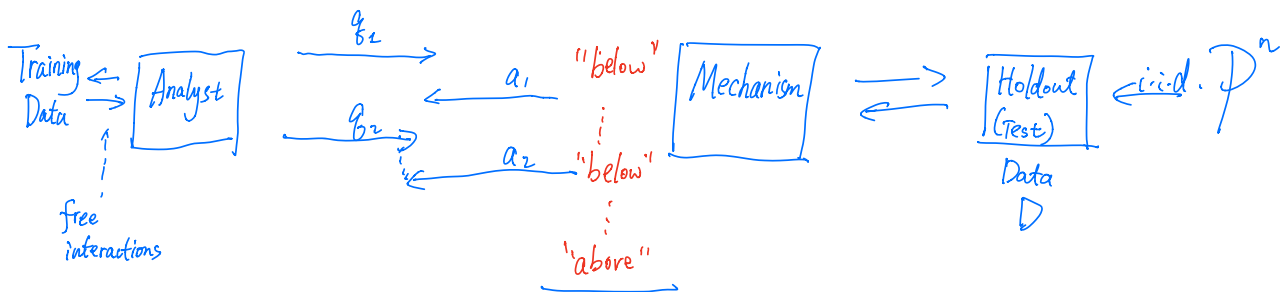
$\rightarrow$  Output Compression. (Sparse Vector)

$\rightarrow$  SV + Private Multiplicative Weights  
 (Synthetic Data Method)

# Model for ADA



$b_j$  can depend on  
 $(b_1, a_1, b_2, a_2, \dots, b_{j-1}, a_{j-1})$ .



Example: 
$$g_1(D) = \frac{1}{n} \sum_{i=1}^n [\mathbb{I}[h(x_i) = y_i]]$$

Reusable Holdout.

# Sparse Vector

Input: an adaptive sequence  $f_1, f_2, \dots$   
( $\Delta$ -sensitive)

Dataset  $D$ . Threshold  $T$ .

Noisy Threshold:  $\tilde{T} = T + Z_0, Z_0 \sim \text{Lap}(\frac{2\Delta}{\epsilon})$

For every query  $i$ :

$$\tilde{a}_i = f_i(D) + Z_i, \quad Z_i \sim \text{Lap}(\frac{4\Delta}{\epsilon})$$

If  $\tilde{a}_i < \tilde{T}$ , return  $b_j = \perp$  "Below"

Else return  $b_j = \top$  ; Break.  
"Above" ; Stop.

Binary coding:  $\vec{b} = (\underbrace{000 \dots 0}_{\text{"Below"}} \underbrace{1}_{\text{"Above"}})$

# Accuracy.

Given  $f_1, f_2, \dots$

Return  $b_1, b_2, \dots$

"Empirical Error"

If  $b_i = \text{"Below"}$ ,  $\max\left(0, \underbrace{f_i(D) - T}_{\text{Want to be small}}\right)$ .

If  $b_i = \text{"Above"}$ ,  $\max\left(0, \underbrace{T - f_i(D)}_{\text{Want small.}}\right)$ .

Theorem. With probability  $1 - \beta$ , when run over a sequence of  $k$  queries.  $SV$  has empirical error

$$\alpha \leq \frac{6 \cdot \Delta \ln\left(\frac{k+1}{\beta}\right)}{\epsilon}$$

Error scales logarithmically in  $k$ .

Proof Sketch. Union Bound Laplace Noise  $|Z_i|$ .

---

Composing  $m$  runs of  $SV$ .

Population error.  $\max\left(0, \underbrace{f_i(P) - T}_{\text{Want to be small}}\right)$  for "Below"

$\max\left(0, \underbrace{T - f_i(P)}_{\text{Want small.}}\right)$  for "Above"



$$\alpha' \leq O\left(\underbrace{\frac{\ln k}{n\varepsilon}} + \underbrace{\varepsilon\sqrt{m}}\right)$$

Transfer  
Theorem

: Empirical error  
for answering  
 $k$  queries in total

Generalization error  
for  $(\varepsilon\sqrt{m}, \delta)$ -DP

$$\leq O\left(\frac{m^{\frac{1}{4}} \ln^{\frac{1}{2}} k}{\sqrt{n}}\right).$$

$m$ : # "Above"

← "Interesting Events"

$k$ :  $\approx$  # "Below"

Privacy.

Theorem. SparseVector is  $(\epsilon, \delta)$ -D.P.

Proof. Sketch.

Fix any neighbors  $D$  &  $D'$ ,  
any output  $\vec{b} = (\perp)^{(k-1)} (\top)$

Given the output,  $q_1, \dots, q_k$  are also fixed.

Furthermore, fix noise values

$$Z_1 = z_1, Z_2 = z_2, \dots, Z_{k-1} = z_{k-1}.$$

$$\text{Let } g(D) = \max_{j=1}^{k-1} q_j(D) + z_j$$

The output is  $\vec{b} = (\perp)^{(k-1)} (\top)$  if and only if

$$g(D) < \tilde{T} \leq q_k(D) + z_k \quad \text{for } D$$

$$\text{and } g(D') < \tilde{T} \leq q_k(D') + z_k \quad \text{for } D'$$

The only randomness we consider:

$$\left( \tilde{T}, Z_k \right).$$

For every realization  $(\tilde{r}, \tilde{z})$ ,

we construct  $(\tilde{r}', \tilde{z}')$  :

$$\tilde{r}' = \tilde{r} + g(D') - g(D)$$

$$\tilde{z}' = \tilde{z} + g(D') - g(D) + g_{\delta k}(D) - g_{\delta k}(D')$$

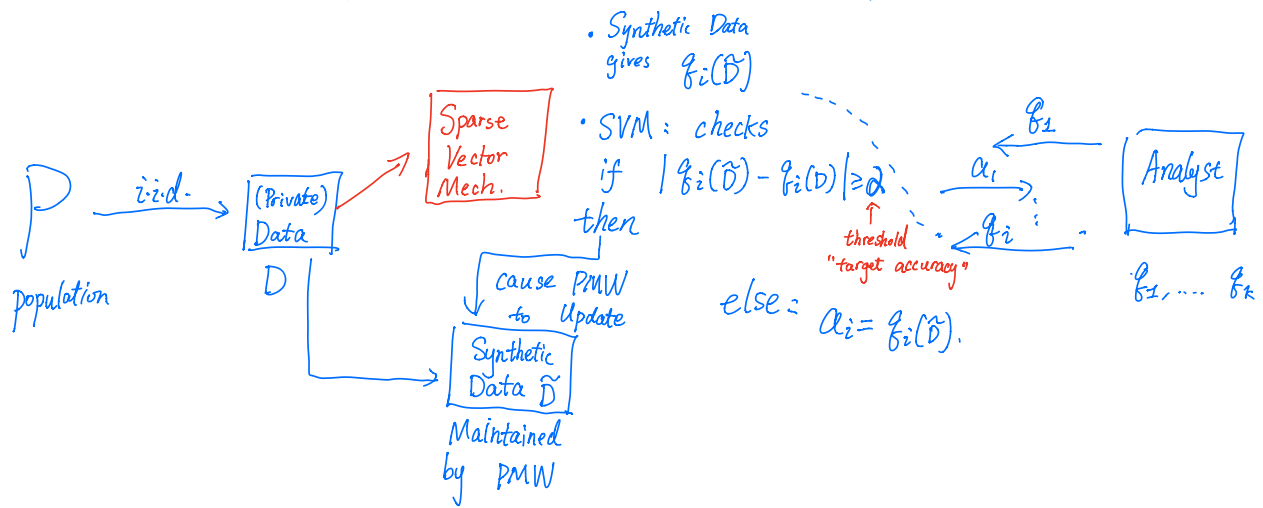
$$P\left[\left(\frac{\tilde{r}}{T}, \tilde{z}_k\right) = (\tilde{r}, \tilde{z})\right] \leq e^\varepsilon P\left[\left(\frac{\tilde{r}}{T}, \tilde{z}_k\right) = (\tilde{r}', \tilde{z}')\right]$$

$$\Downarrow \\ M(D) = \vec{b}$$

$$\Downarrow \\ M(D') = \vec{b}.$$



# Private Multiplicative Weights (PMW)



Hope: Answer most queries using the synthetic data.

PMW: adaptive analog of MWEM  
 (MW w/ Exp Mech)

PMW: runs in time linearly in  $|\mathcal{X}|$ .

$$(\mathcal{X} = \{0,1\}^d)$$

# Sample Complexity / Accuracy.

Non-adaptive Queries.

- Take empirical averages:  $a_j = g_j(D)$   
$$\max_j |a_j - g_j(P)| \approx \sqrt{\frac{\log(k)}{n}}$$

Adaptive Queries

- Sample Splitting Method:  $D_1, \dots, D_k$ ,  $a_j = g(D_j)$   
$$\max_j |a_j - g_j(P)| \leq \sqrt{\frac{k}{n}}$$

Use  
Transfer  
Theorem

- Gaussian Mechanism  
$$\max_j |a_j - g_j(P)| \leq \sqrt{\frac{k^{\frac{1}{2}}}{n}}$$

- Private Multiplicative Weights + Sparse Vector

$$\approx \min \left\{ \left( \frac{\ln(k) \sqrt{\ln|\mathcal{X}|}}{n} \right)^{\frac{1}{3}}, \sqrt{\frac{k^{\frac{1}{2}}}{n}} \right\}$$

Remark: • Runtime is Linear in  $|\mathcal{X}|$  (Infinite or  $\exp(d)$ )

- The error bound depends on  $\log|\mathcal{X}| \approx d$ .  
↑  
dimension of data.

Lower Bound. Statement

" If either the algorithm is polytime (d)  
efficient

or  $\ln|\mathcal{X}| \geq \Omega(n^2)$ ,

then # queries  $\underline{k} \leq \mathcal{O}(n^2)$ .

(before analyst comes up w/ a high-error query)

Satisfies.

Gaussian Mech  $k \approx n^2$ .

11