# Lecture 23

## Adaptive Data Analysis.

Method → Sample → Conclusion → Revise Method / hypothesis (loops back to Method)
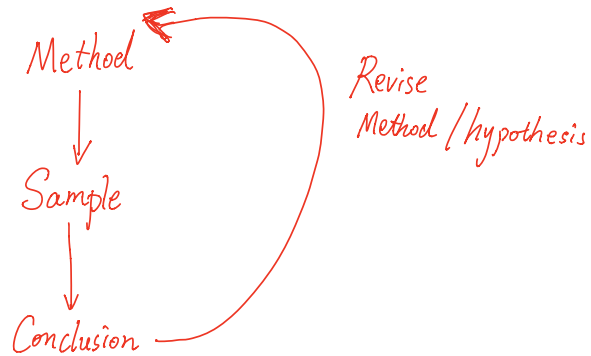
Logistics:
Project Presentation   (May 3 & 5)
            Schedule  Announced  this  week
        20   mins.
No   more   HW.

# Model for Adaptive Data Analysis

## Statistical / Linear Queries

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

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

$\uparrow$ population

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

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

Example: mean, correlation, variance, error/risk, gradients

Extension:

① Low-Sensitive Queries

$$\forall \text{ neighbors } D \& D', \qquad |q(D) - q(D')| \leq \Delta$$

② Minimization Queries

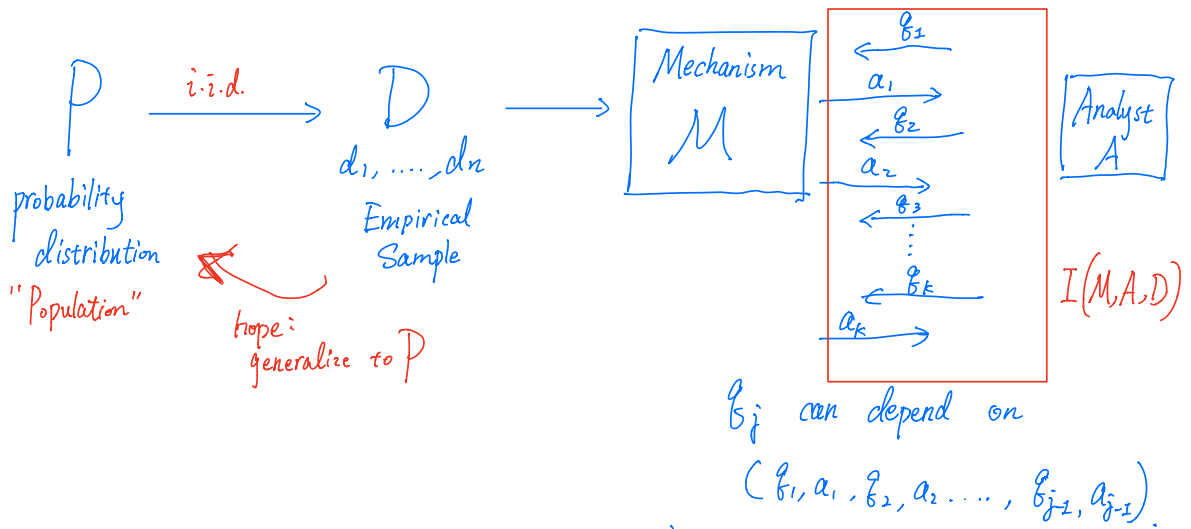Query is given by some loss function $L : \mathcal{X}^n \times \Theta \longrightarrow [0,1]$

$\uparrow$ Dataset $\quad \llcorner$ Parameter Space.

Answer $\theta \in \Theta$

$$\forall \text{ neighbors } D \& D' \qquad |L(D,\theta) - L(D',\theta)| \leq \Delta.$$

$$\forall \theta \in \Theta$$

# Interaction of Adaptive Data Analysis.

$$P \xrightarrow{i.i.d.} D \longrightarrow$$

probability
distribution
"Population"

$d_1, \ldots, d_n$
Empirical
Sample

hope:
generalize to $P$

Mechanism
$M$

$$\begin{array}{l} \overleftarrow{\;q_1\;} \\ \xrightarrow{\;a_1\;} \\ \overleftarrow{\;q_2\;} \\ \xrightarrow{\;a_2\;} \\ \overleftarrow{\;q_3\;} \\ \vdots \\ \overleftarrow{\;q_k\;} \\ \xrightarrow{\;a_k\;} \end{array}$$

Analyst
$A$

$I(M,A,D)$

$q_j$ can depend on

$$( q_1, a_1, q_2, a_2 \ldots, q_{j-1}, a_{j-1} ).$$

Transcript $\Pi = \big( (q_1, a_1), \ldots, (q_k, a_k) \big) \longleftarrow I(M A, D)$

$\uparrow$
Interaction

"Goal" : $\forall j$ ,

$$\big| a_j - q_j(P) \big| \leq \text{small}.$$

$\uparrow$ Population Value.

Not just empirical averages

Avoid : Queries $q$ s.t.

$$\big| q(P) - q(D) \big| \geq \text{Large}$$

$\uparrow$

$D$ is not representative.

# DP $\implies$ Generalization in ADA

$(\alpha, \beta)$ – sample accuracy

$$\Pr_{D \sim p^n, \, \Pi} \left[ \max_j \left| q_{f_j}(D) - a_j \right| \geq \alpha \right] \leq \beta$$

$(\alpha, \beta)$ – distributional accuracy

$$\Pr_{D \sim p^n, \, \Pi} \left[ \max_j \left| q_{f_j}(P) - a_j \right| \geq \alpha \right] \leq \beta$$

# Sample Complexity / Accuracy.

### Non-adaptive Queries.

- Take empirical averages: $a_j = q_j(D)$

$$\max_j \left| a_j - q_j(P) \right| \lesssim \sqrt{\frac{\log(k)}{n}}$$

### Adaptive Queries

- Sample Splitting Method: $D_1, \ldots, D_k$, $a_j = q(D_j)$

$$\max_j \left| a_j - q_j(P) \right| \lesssim \sqrt{\frac{k}{n}}$$

- Differential Privacy. Gaussian Mechanism

$$a_j = q_j(D) + N(0, \sigma^2).$$

$$\max_j \left| a_j - q_j(P) \right| \lesssim \frac{k^{\frac{1}{4}}}{\sqrt{n}} \qquad O(\alpha + \varepsilon).$$

$$\underbrace{\frac{\sqrt{k}}{n\sigma}}_{\substack{\varepsilon \\ \text{privacy/Gen} \\ \text{Bound}}} + \underbrace{\sigma}_{\substack{\alpha \\ \text{sample} \\ \text{Accuracy Bound.}}}$$

Adding Noise
Reduces Error.

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

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

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

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

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

---

Simpler version $(\varepsilon, 0)$-DP

$(\alpha, \beta)$ - sample accuracy
$(\varepsilon, 0)$ - DP.

$\Rightarrow (\alpha', \beta')$ - distributionally accurate.

$$\alpha' = \alpha + (e^\xi - 1) + \sqrt{\frac{2\ln(1/\eta)}{n}} \quad , \quad \beta' = \beta + \eta \, , \text{ for all } \eta > 0.$$

$$\approx \alpha + \varepsilon + \hat{O}\left(\frac{1}{\sqrt{n}}\right)$$

$\underset{\text{Sampling error bound}}{\upharpoonleft}$

# Proof Sketch.

$$D \xleftarrow{\;i.i.d.\;} P^n$$

Transcript: $\pi = (q_1, a_1, \ldots, q_k, a_k) \longleftarrow I(M, A, D)$

$$Q_\pi = P^n \big/ \pi$$

" posterior distribution over $D$ conditioned on $\pi$"

- Suppose $M$ is $(\alpha, \beta) -$ sample accurate.

$$\mathbb{P}\left[\; \underbrace{\max_{j} |a_j - q_j(D)| \geq \alpha}\; \right] \leq \beta$$

$\qquad\qquad$ Event $E$ about $(D, \pi)$

## Lemma ( Bayesian Resampling )

$$\mathop{\mathbb{P}}_{\substack{D \sim P^n \\ \pi \leftarrow I(M, A, D)}}\left[ (D, \pi) \in E \right] = \mathop{\mathbb{P}}_{\substack{D \sim P^n \\ \pi \sim I(M, A, D) \\ D' \sim Q_\pi}}\left[ (D', \pi) \in E \right]$$

- Generic : Nothing to do w/ DP.

- Sample accuracy w.r.t. $D$

$$\Longrightarrow \text{ Sample accuracy w.r.t. } D' \sim Q_\pi$$

- $q(D) \approx q(Q_\pi) = \mathop{\mathbb{E}}_{D' \sim Q_\pi}\left[ q(D') \right]$

- Goal: $\quad q(D) \approx q(P)$
$$\qquad\qquad\; a \approx q(P)$$

Missing Step: $\quad q(Q_\pi) \approx q(P)$

Do you want to see the proof?

Proof. $\mathbb{P}_{\substack{D \sim p^n \\ \Pi \leftarrow I(M,A,D) \\ D' \sim Q_\pi}} \left[ (D', \pi) \in E \right]$

$$= \sum_{D=x} \sum_{\Pi=\pi} \sum_{D'=x'} \mathbb{1}\left[(x',\pi) \in E\right] \cdot \mathbb{P}\left[D=x\right] \cdot \mathbb{P}\left[\pi \mid x\right] \cdot \mathbb{P}\left[x' \mid \pi\right]$$

$$= \sum_{D=x} \sum_{\Pi=\pi} \sum_{D'=x'} \mathbb{1}\left[(x',\pi) \in E\right] \cdot \mathbb{P}\left[D=x, \Pi=\pi\right] \underset{\substack{D' \sim p^n}}{\mathbb{P}}\left[D'=x' \mid \Pi=\pi\right]$$

$$= \sum_{\Pi=\pi} \sum_{D'=x'} \mathbb{1}\left[(x',\pi) \in E\right] \cdot \mathbb{P}\left[\Pi=\pi\right] \cdot \underset{\substack{D' \sim p^n \\ \Pi \leftarrow I}}{\mathbb{P}}\left[D'=x' \mid \Pi=\pi\right]$$

$$= \sum_{\Pi=\pi} \sum_{D'=x'} \mathbb{1}\left[(x',\pi) \in E\right] \cdot \mathbb{P}\left[D'=x'\right] \cdot \underset{\substack{D' \sim p^n \\ \Pi \leftarrow I}}{\mathbb{P}}\left[\Pi=\pi \mid D'=x'\right]$$

$$= \underset{\substack{D \sim p^n \\ \Pi \leftarrow I(M,A,D)}}{\mathbb{P}}\left[(D,\Pi) \in E\right]$$

Differential Privacy $\implies$ $f(Q_\pi) \approx f(P)$

Proof (Sketch):

$$f(Q_\pi) = \mathbb{E}_{D' \sim Q_\pi} \left[ f(D') \right]$$

$$= \mathbb{E}_{\substack{D' \sim Q_\pi \\ i \leftarrow \text{unif}\{1,\dots n\}}} \left[ f(d'_i) \right] = \mathbb{E}_{\substack{D \sim P^n \\ i \leftarrow \text{unif}\{1 \dots n\}}} \left[ f(d_i) \right]$$

$$= \int_{x \in \mathcal{X}} f(x) \cdot \mathbb{P}_{\substack{D \sim P^n \\ i \sim \text{unif}\{1,\dots n\}}} \left[ d_i = x \mid \pi \right] \quad \xleftarrow{\text{Conditioning}}$$

Bayes Rule

$$= \int_{x \in \mathcal{X}} f(x) \cdot \frac{\mathbb{P}_{\substack{D \sim P^n \\ i \leftarrow \text{unif}\{1,\dots n\}}} \left[ \pi \mid d_i = x \right] \cdot \mathbb{P}_{\substack{D \sim P^n \\ i \leftarrow \{1,n\}}} \left[ d_i = x \right]}{\mathbb{P}_{D \sim P^n} \left[ \pi \right]}$$

Approximate Cancellation?

$\pi \leftarrow I(M, A, D)$
$\uparrow$
$(\varepsilon, 0)\text{-DP}$

$$\leq \int_{x \in \mathcal{X}} f(x) \; \frac{e^\varepsilon \, \mathbb{P}[\pi] \cdot \mathbb{P}[d_i = x]}{\mathbb{P}[\pi]}$$

$$= e^\varepsilon \cdot \int_{x \in \mathcal{X}} f(x) \cdot \mathbb{P}_{\substack{D \sim P^n \\ i \sim \{1,\dots n\}}} \left[ d_i = x \right] = f(P).$$

$$f(Q_\pi) \leq e^\varepsilon \; f(P)$$

$$f(Q_\pi) \geq e^{-\varepsilon} \; f(P).$$

# Putting Together

- $f(Q_\pi) \in \left[ e^{-\varepsilon} f(P), \ e^{\varepsilon} f(P) \right]$

  $\nearrow$
  Due to $(\varepsilon, 0) - DP.$

- $" f(D) \approx f(Q_\pi)"$

$$\Pr_{\substack{D \leftarrow P^n \\ \pi \leftarrow I}} \left[ |f(Q_\pi) - f(D)| \geq C_\eta \right] = \Pr_{\substack{D \leftarrow P^n \\ \pi \leftarrow I \\ D' \leftarrow Q_\pi}} \left[ |f(Q_\pi) - f(D')| \geq C_\eta \right] \leq \eta$$
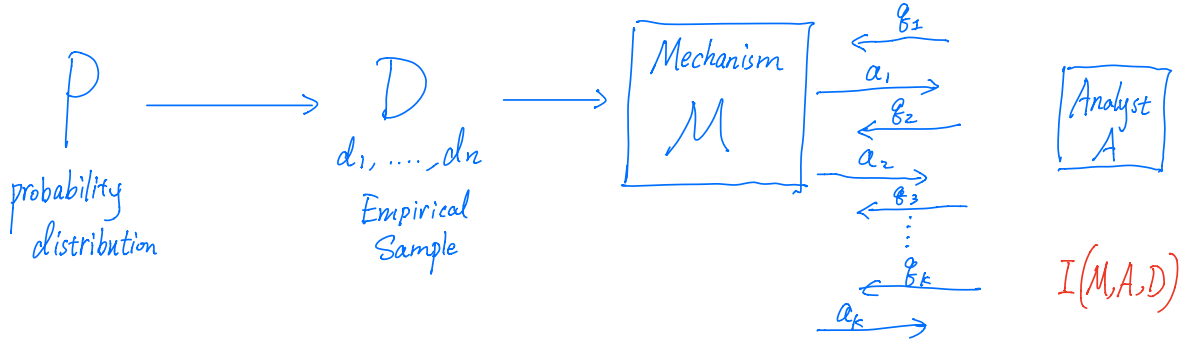
$$C_\eta = \sqrt{\frac{2 \ln(1/\eta)}{n}} \qquad \substack{\text{Azuma} \\ \text{Ineq.}}$$

- Sample accuracy, w.p. $1 - \beta.$

  $|a - f(D)| \leq \alpha$ for all queries. $\beta = \beta_{ij}.$

# Add 3 sources of error $\boxed{\ }$

# Mental Model of ADA.

$$P \longrightarrow D \longrightarrow \boxed{\begin{array}{c} \text{Mechanism} \\ M \end{array}}$$

probability
distribution

$d_1, \ldots, d_n$

Empirical
Sample

$\boxed{\begin{array}{c} \text{Analyst} \\ A \end{array}}$

$\xleftarrow{q_1}$
$\xrightarrow{a_1}$
$\xleftarrow{q_2}$
$\xrightarrow{a_2}$
$\xleftarrow{q_3}$
$\vdots$
$\xleftarrow{q_k}$
$\xrightarrow{a_k}$

$I(M, A, D)$

$q_j$ can depend on

$(q_1, a_1, q_2, a_2 \ldots, q_{j-1}, a_{j-1})$.

---

Training
Data $\xrightleftharpoons{}$ $\boxed{\text{Analyst}}$

$\xrightarrow{q_1}$

$\xleftarrow{a_1}$ "below"

$\xrightarrow{q_2}$

$\xleftarrow{a_2}$ "below"

"above"

$\boxed{\text{Mechanism}}$ $\xrightarrow{}$ $\xleftarrow{}$ $\boxed{\begin{array}{c} \text{Holdout} \\ \text{(Test)} \end{array}}$ $\xleftarrow{i.i.d.} P^n$

Data
$D$

free
interactions

Example: $\quad q_1(D) = \dfrac{1}{n} \displaystyle\sum_{i=1}^{h} \left[ \mathbb{1}\left[ h(x_i) = y_i \right] \right]$

error
accuracy

## Reusable Holdout.

# Sparse Vector   or   Thresholdout

---

**Algorithm 2: SparseVector$(\mathbf{s}, T, \Delta, \epsilon, q_1, q_2, \ldots)$:**

---

**Input:** $q_1, q_2 \ldots$ is a stream of $\Delta$-sensitive queries

1 **AllDone** $\leftarrow$ **FALSE**;     $T$: threshold     Stop when

2 $\tilde{T} = T + Z_0$ where $Z_0 \sim \mathrm{Lap}(2\Delta/\epsilon))$ ;     $\frac{q(s)}{b}$ is above $T$. "$\frac{q}{b}(D)$"

3 **while** *not* **AllDone do**

4      Accept the next query $q_i$;

5      $a_i \leftarrow q_i(\mathbf{s})$ ;

6      $\tilde{a}_i \leftarrow a_i + Z_i$ where $Z_i \sim \mathrm{Lap}(4\Delta/\epsilon)$ ;

7      **if** $\tilde{a}_i < \tilde{T}$ **then**

8          **return** $b_j = \bot$;

9      **else**

10          **return** $b_j = \top$ ;

11          **AllDone** $\leftarrow$ **TRUE** ;

---

- Only Noise $\top$ once.
- Could release "$\bot$" many times.