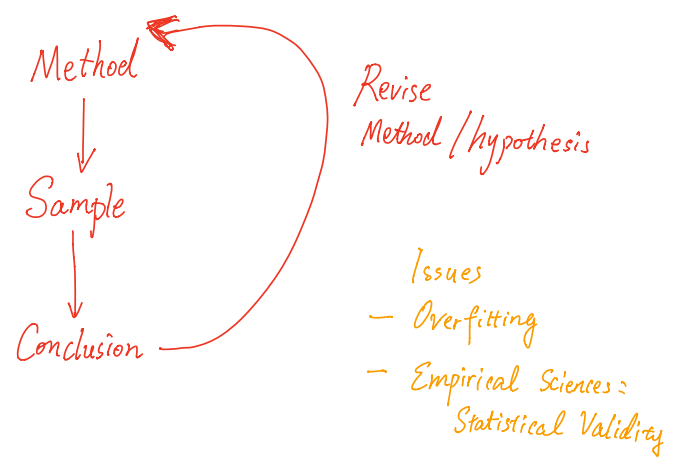


Lecture 22

Adaptive Data Analysis.



Case Study: "Wacky Boost"

Suppose $\{x, y\} \in \{0, 1\}^d \times \{0, 1\}$
 $\uparrow \qquad \qquad \qquad \uparrow$
d-dim features Binary labels

Learn $f: X \mapsto Y$

Accuracy: $Acc(f) = \mathbb{P}_{(x,y) \sim P} [f(x) = y]$ ← proxy?

Dataset $D: Acc_D(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[f(x_i) = y_i]$

ALG WB:

$F = \{ \}$ "selected feature"

For $j = 1, \dots, d$

Compute $C_j = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[x_j^{(i)} = y^{(i)}]$

If $C_j \geq \frac{1}{2} + \frac{1}{\sqrt{n}}$, then $F \leftarrow F \cup \{j\}$

ENDFOR

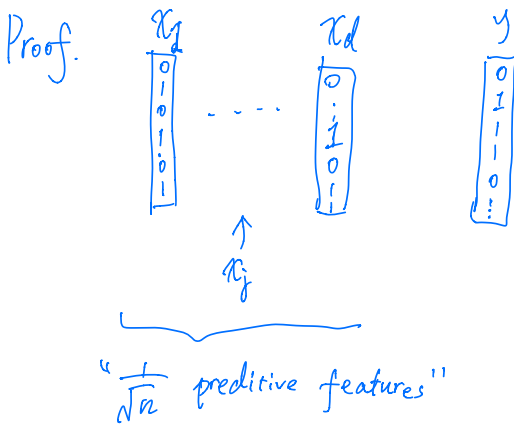
Output: $\hat{f}(x) = \begin{cases} 1 & \text{if } \sum_{j \in F} x_j \geq \frac{|F|}{2} \\ 0 & \text{o/w.} \end{cases}$ "Majority Vote"

Theorem. Let \mathcal{P} denote the uniform distribution over $\{0, 1\}^d \times \{0, 1\}$. There exists a constant c such that with probability $1 - \delta$,

$$|\text{Acc}_D(\hat{f}) - \text{Acc}(\hat{f})| \geq 0.49$$

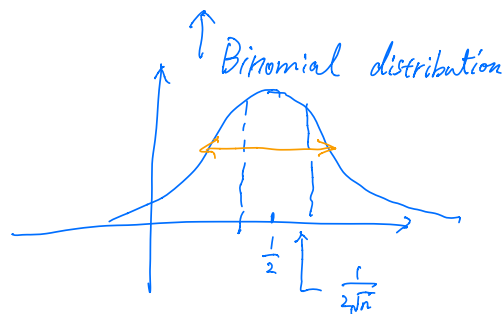
as long as $d \geq c \cdot \max(n, \log(\frac{1}{\delta}))$.

- None of the features is predictive.
- $\text{Acc}(\hat{f}) \approx \frac{1}{2}$, $\text{Acc}_D(\hat{f}) \rightarrow 99\%$. \rightarrow gap of 49%.



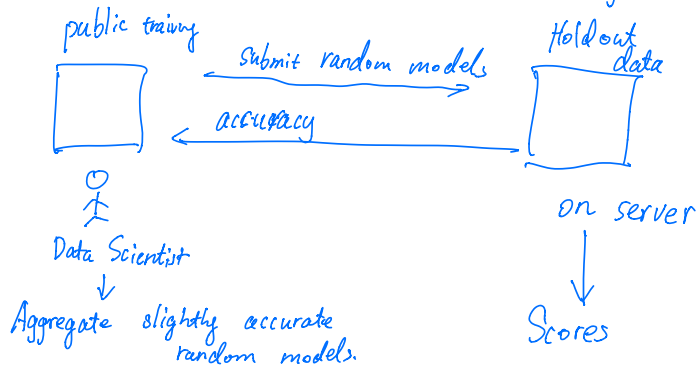
\rightarrow 99% predictor.

$$C_j = \frac{1}{n} \sum_i \mathbb{1}[x_j^{(i)} = y^{(i)}]$$



$C_j \geq \frac{1}{2} + \frac{1}{\sqrt{n}}$ happens with $\Omega(1)$ prob.

Background: Kaggle.



Model for Adaptive Data Analysis

Statistical / Linear Queries

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

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

↑
population

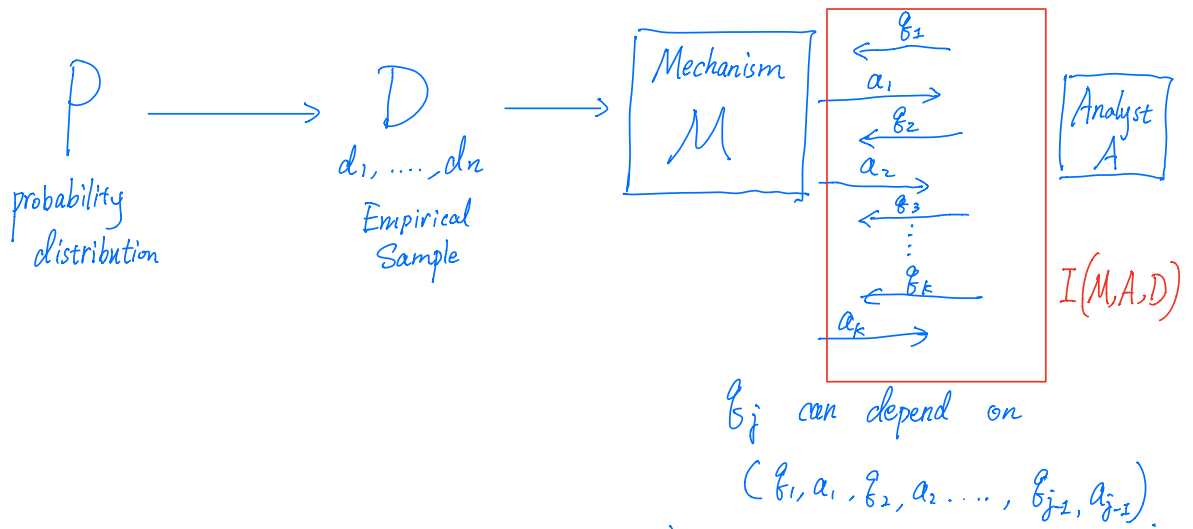
$$\hat{\mathbb{E}}_{d_i \sim D} [\phi(d_i)] = \frac{1}{n} \sum_{i=1}^n [\phi(d_i)] \quad \text{"Empirical Average"}$$

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

Why SQ's ?

- Mean, Variances, correlations, etc.
- Error / Risk of Predictive Models
$$\mathbb{E}[\ell(f(x), y)]$$
- Gradient of loss of a hypothesis
$$\mathbb{E}[\nabla \ell(f(x), y)]$$
- Statistical Query Model (Kearns '94).
"pretty much" all PAC learning problems

Interaction of Adaptive Data Analysis.



Transcript $\Pi = (g_1, a_1, \dots, g_k, a_k) \leftarrow I(M, D)$

"Goal" : $\forall j,$

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

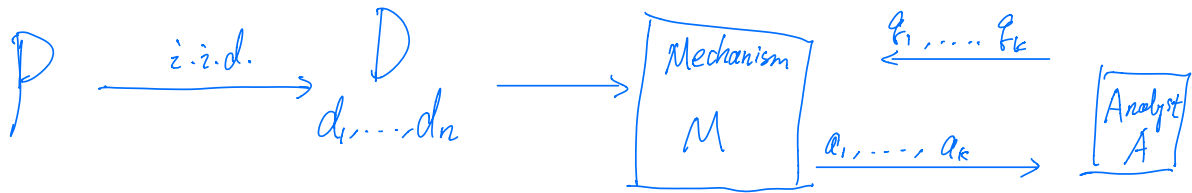
\uparrow Population Value.
Not just empirical averages

Avoid : queries g s.t.

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

\uparrow
D is not representative.

Non-adaptive queries



What would M be?

Output $f(D)$

Theorem. $\max_{j \in \{1, \dots, k\}} |f_j(D) - f_j(P)| \leq \sqrt{\frac{\ln(\frac{2k}{\delta})}{2n}}$ w.p. $1 - \delta$.

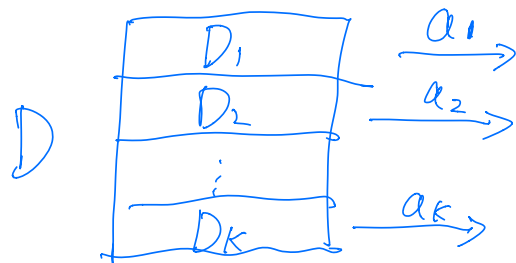
Proof Sketch. Chernoff Bound, $\forall j$

$$P\left[|f_j(D) - f_j(P)| > \sqrt{\frac{\ln(\frac{2k}{\delta})}{n}}\right] \leq \frac{\delta}{k}.$$

Union Bound

\Rightarrow Stated Claim.

Data Splitting



Independence.

$$|D_j| = \frac{n}{k}$$

$$\max_j |f_j(D) - f_j(P)| \leq \sqrt{\frac{k \ln(k)}{n}}$$

DP \implies Generalization in ADA

(α, β) - sample accuracy

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

\uparrow
observe

(α, β) - distributional accuracy

Ultimate Goal.

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

\uparrow
unobserve

Idea: Make sure D is representative w.r.t. f

$$f(D) \approx f(P)$$

\uparrow
Differential Privacy

Transfer Theorem. (ϵ, δ) -version [JLNRSS20]

Suppose $I(M, A, D)$ is (α, β) -sample accurate \leftarrow
& (ϵ, δ) -DP. \leftarrow

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

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

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

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

(α, β) -sample accuracy

$(\epsilon, 0)$ -DP.

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

$$\alpha' = \alpha + (e^\epsilon - 1) + \sqrt{\frac{2 \ln(1/\eta)}{n}}, \quad \beta' = \beta + \eta.$$

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

Sample Complexity / Accuracy.

Non-adaptive Queries.

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

$$\max_j |a_j - f_j(P)| \approx \sqrt{\frac{\log(k)}{n}}$$

Adaptive Queries

- Sample Splitting Method: D_1, \dots, D_k , $a_j = f(D_j)$

$$\max_j |a_j - f_j(P)| \leq \sqrt{\frac{k \log(k)}{n}}$$

- Differential Privacy, Gaussian Mechanism

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

$$\max_j |a_j - f_j(P)| \approx \underbrace{\frac{k^{\frac{1}{4}}}{\sqrt{n}}}_{\epsilon} \quad O(\alpha + \epsilon).$$

Adding Noise

Decreases Error!

$$\frac{\sqrt{k}}{n\sigma} + \sigma$$

ϵ privacy/gen Bound

α sample Accuracy Bound.

