# BU CS 591 Spring 2025
# Privacy in ML and Statistics

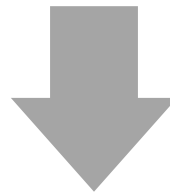**Adam Smith (BU)**

**Jonathan Ullman (NEU)**

**Lecture 24: Adaptive Data Analysis**

# Today

- ## Adaptive validity in statistical analysis
  - Example setting: ML competitions
  - What can go wrong
  - Nothing about privacy!

- ## Privacy prevents overfitting
  - Single query case
  - Extension to multiple queries
  - General transfer theorem

# Statistical Theory

Method

⬇

Sample (from population)

⬇

Conclusions

Statistical analysis guarantees that your conclusions generalize to the population

# Statistical Practice

## Why Most Published Research Findings Are False

John P. A. Ioannidis

# The Statistical Crisis in Science

Data-dependent analysis—a "garden of forking paths"— explains why many statistically significant comparisons don't hold up.

Andrew Gelman and Eric Loken

# Statistical Practice



Method

Sample

Conclusions

Statistical guarantees no longer apply
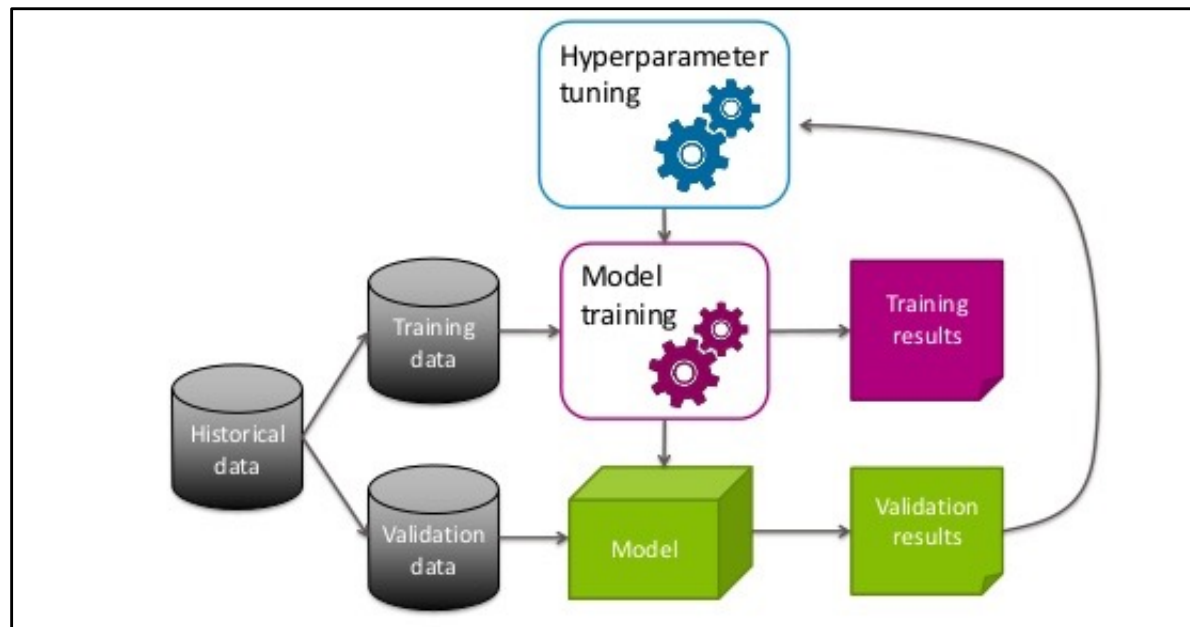when the method and sample are correlated

# Examples of Adaptive Data Analysis

Well-specified adaptive algorithms
   Select features then fit a model (Freedman's Paradox)
   Hyperparameter tuning (sometimes)
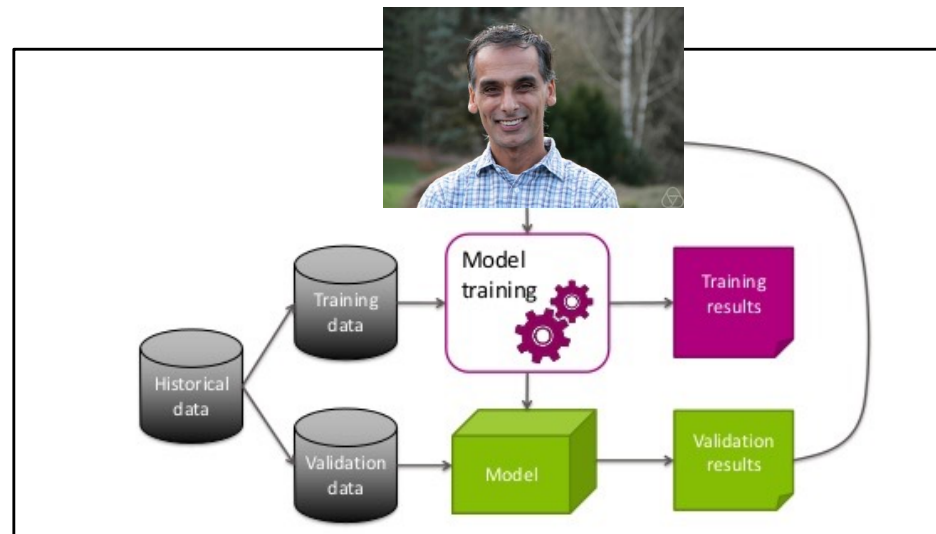   **Data science competitions**



Alice Zheng.  "Evaluating Machine Learning Models."

# Examples of Adaptive Data Analysis
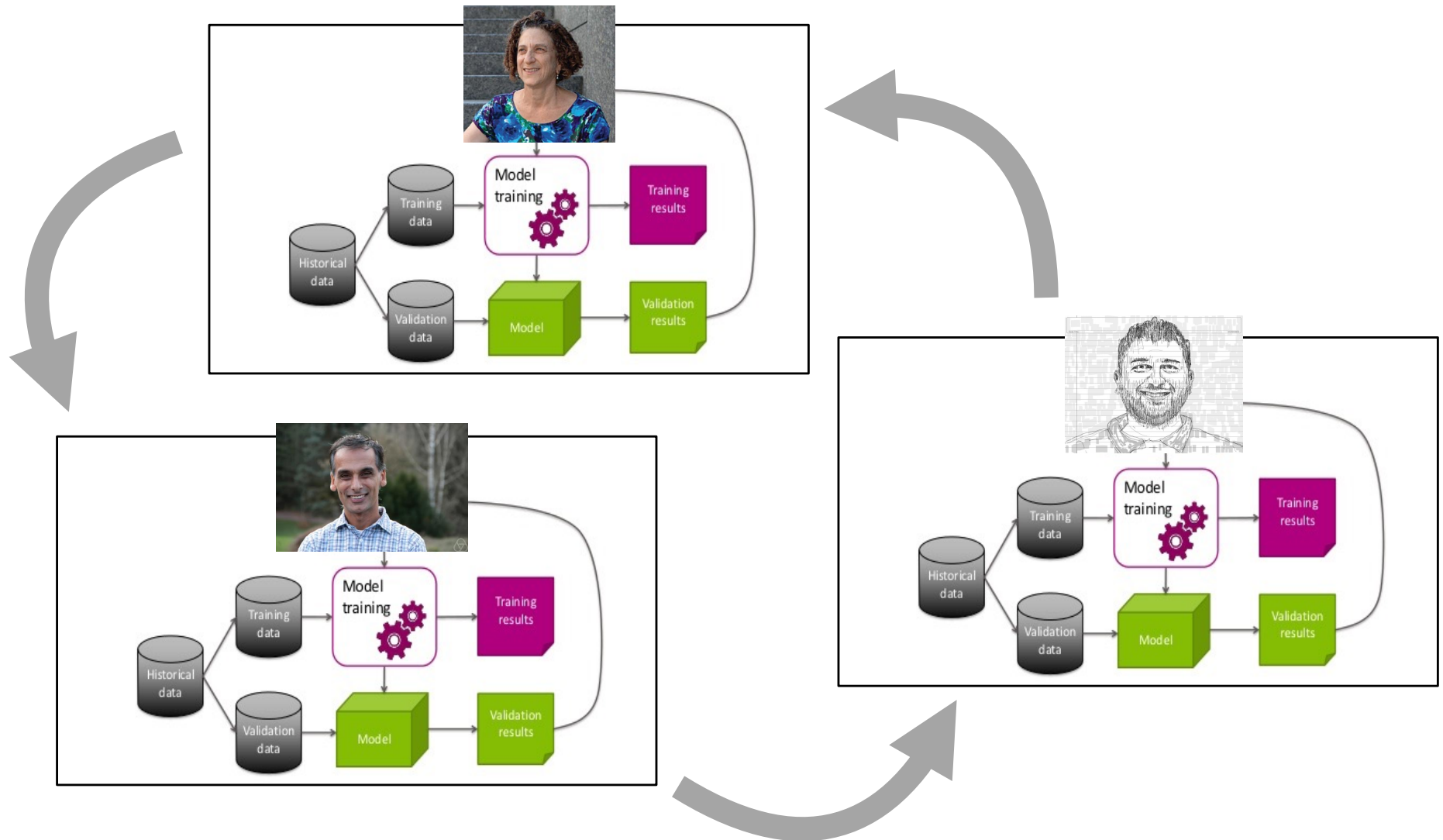
## Researcher degrees of freedom

> The interaction effect is not significant when the scale from the Danish study are used to gauge the US subjects' support for redistribution. This arises because two of the items are somewhat unreliable in a US context. Hence, for items 5 and 6, the inter-item correlations range from as low as .11 to .30. These two items are also those that express the idea of European-style market intervention most clearly and, hence, could sound odd and unfamiliar to the US subjects. When these two unreliable items are removed ($\alpha$ after removal = .72), the interaction effect becomes significant.

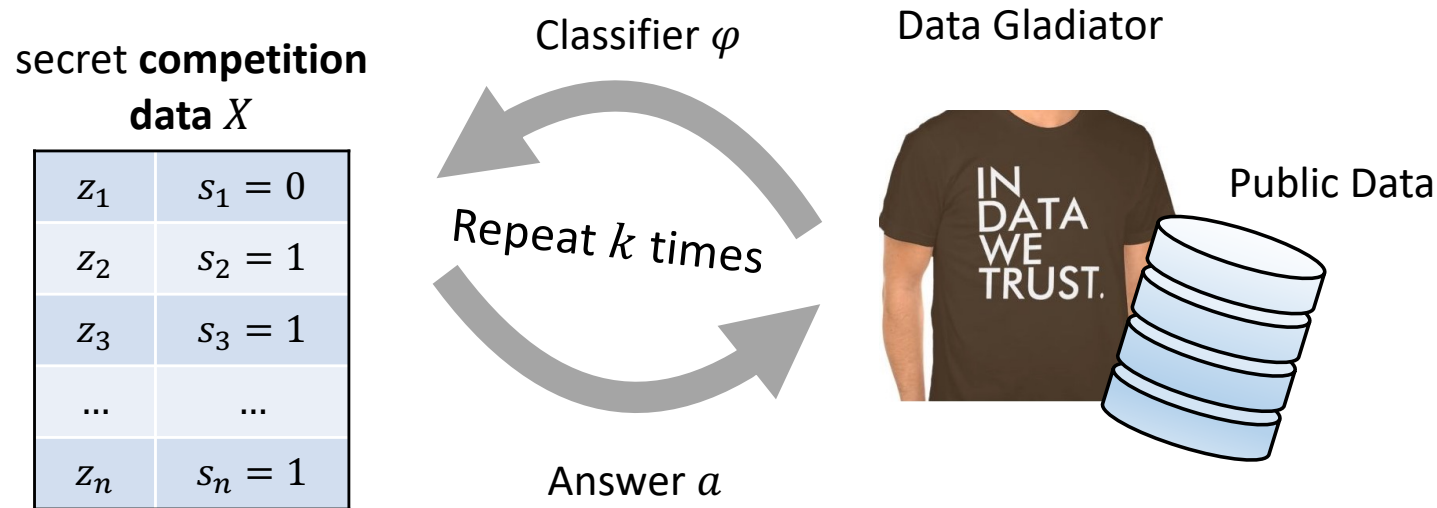A. Gelman, E. Loken. "The Garden of Forking Paths."

# Examples of Adaptive Data Analysis

Reuse of datasets by multiple researchers

# Case Study: ML Competitions

**kaggle**

secret **competition data** $X$

| | |
|---|---|
| $z_1$ | $s_1 = 0$ |
| $z_2$ | $s_2 = 1$ |
| $z_3$ | $s_3 = 1$ |
| ... | ... |
| $z_n$ | $s_n = 1$ |

Classifier $\varphi$

Repeat $k$ times

Data Gladiator

IN DATA WE TRUST.

Public Data

Answer $a$

$$a \approx \mathrm{score}_X(\varphi) = \frac{1}{n}\sum_i \mathbf{1}\{\varphi(z_i) = s_i\} = \mathbb{E}_X(\mathbf{1}\{\varphi(z_i) = s_i\})$$
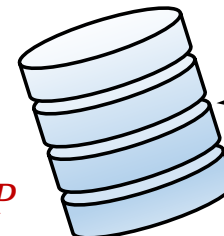
where $\varphi$ is a classifier

Competition: find a classifier $\varphi^*$
with large score **on the distribution**

> Needed: a method for estimating the score
> **on the underlying distribution**

$\mathrm{score}_P(\varphi) = \mathbb{E}_P(\mathbf{1}\{\varphi(z_i) = s_i\})$
score on the underlying
distribution

Secret Prize
Distribution $P$

> Competition
> distribution drawn
> from $P$

# Case Study: ML Competitions

**kaggle**

- Suppose prize and competition data have **random labels**

    - Any classifier will have $\text{score}_P(\varphi) \approx \frac{1}{2}$ on the prize distribution $P$

    - If $\text{score}_X(\varphi) \gg \frac{1}{2}$ then we have overfit

- **How can we prevent the competitors from overfitting to the competition data?**

- **Naïve algorithm:**

    - answer $a = \text{score}_X(\varphi) = \frac{1}{n}\sum_i \mathbf{1}\{\varphi(z_i) = s_i\}$

    - Let's see how well this algorithm does at preventing overfitting

# Non-adaptive analysis

kaggle

- **Competitor's strategy (non-adaptive):**
  - Choose $k$ random classifiers $\varphi_1, \ldots, \varphi_k$
  - Receive $a_1, \ldots, a_k$ where $a_j = score_X(\varphi_j)$
  - Output $\varphi^* = \mathrm{argmax}\, score_X(\varphi_j)$

number of samples (n) = 1000

95% significance threshold

**Theorem (nonadaptive accuracy):**

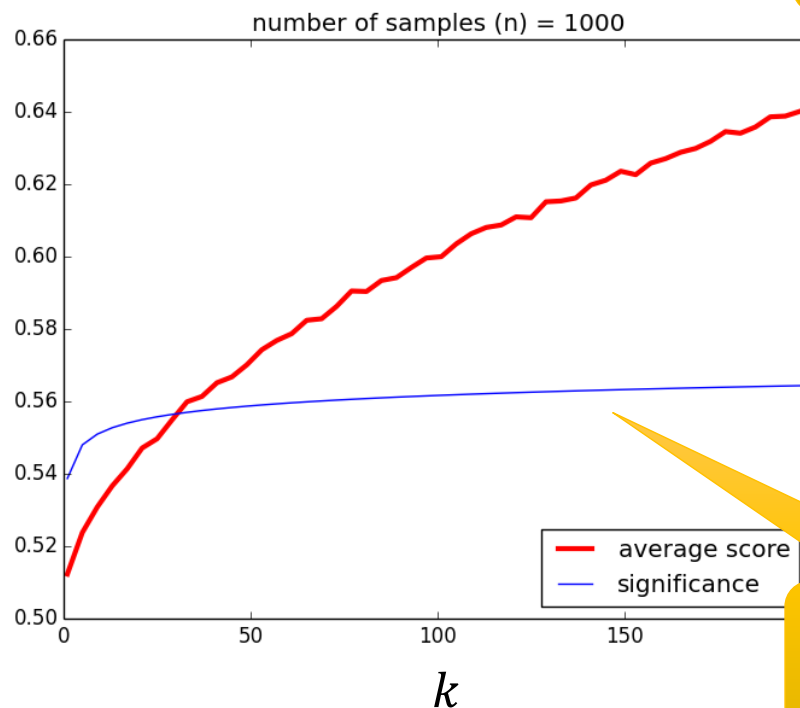$$\mathbb{E}\left(\max_j sc_X(\varphi_j) - sc_P(\varphi_j)\right) \leq \sqrt{\frac{C \cdot \ln k}{n}}$$

1/2

average score

significance

$k$

# Overfitting with adaptive analysis

**kaggle**

- **Competitor's strategy (adaptive):**

  - Choose $k-1$ random classifiers $\varphi_1, \dots, \varphi_{k-1}$
    Receive scores $a_1, \dots, a_{k-1}$

  - Define $\varphi_k(z) = \text{sign}\left(\sum_j \left(a_j - \frac{1}{2}\right) \cdot \varphi_j(z)\right)$

Deviation from population mean

number of samples (n) = 1000

- average score
- significance

$k$

**Theorem (adaptive attack on raw scores):**

$$\mathbb{E}\left(\text{sc}_X(\varphi_k) - \text{sc}_P(\varphi_k)\right) = \Omega\left(\sqrt{\frac{k}{n}}\right)$$

95% significance threshold

# Overfitting with adaptive analysis
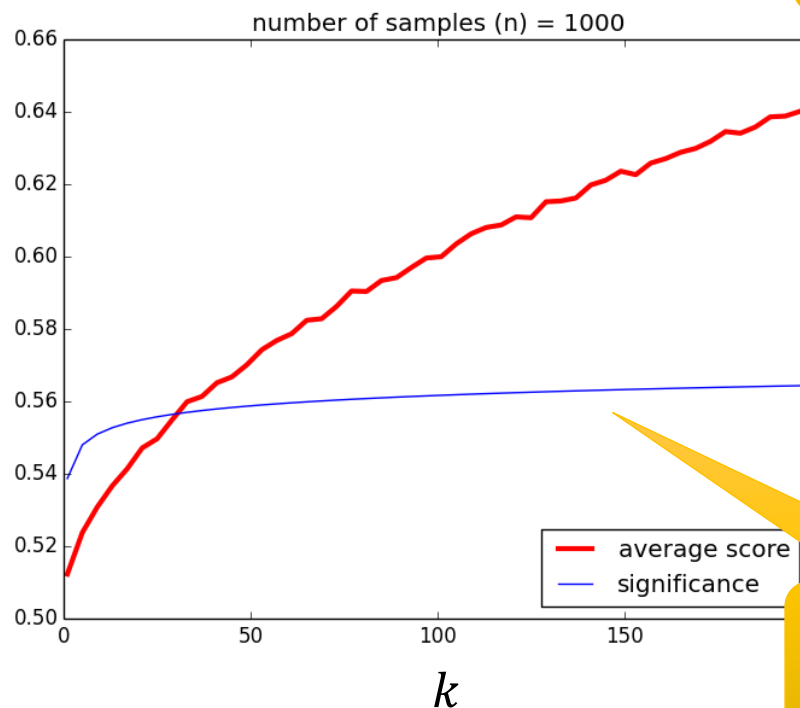
**kaggle**

- **Competitor's strategy (adaptive):**
  - Choose $k-1$ random classifiers $\varphi_1, \ldots, \varphi_{k-1}$
    Receive scores $a_1, \ldots, a_{k-1}$
  - Define $\varphi_k(z) = \text{sign}\left(\sum_j \left(a_j - \frac{1}{2}\right) \cdot \varphi_j(z)\right)$ $= \text{sign}\left\langle \vec{a} - \frac{\vec{1}}{2} , \vec{y} \right\rangle$

*[handwritten annotations:]* $\varphi_k$ "runs" a membership inference attack!

$\vec{y} = \left(\varphi_1(z), \ldots, \varphi_{k-1}(z)\right)$

> Deviation from population mean



number of samples (n) = 1000

Legend:
— average score
— significance

> 95% significance threshold

**Theorem (adaptive attack on raw scores):**

$$\mathbb{E}\left(\text{sc}_X(\varphi_k) - \text{sc}_P(\varphi_k)\right) = \Omega\left(\sqrt{\frac{k}{n}}\right)$$

# What Happened in This Example?

# Case Study: ML Competitions

- **Improved estimator:** Add Gaussian noise $N(0, \sigma^2)$ to the estimated score of each classifier
  - Give answers $a_j = \text{score}_X(\varphi_j) + N(0, \sigma^2)$

# Case Study: ML Competitions

**kaggle**™

- **Improved estimator:** Add Gaussian noise $N(0, \sigma^2)$ to the estimated score of each classifier
  - Give answers $a_j = \text{score}_X(\varphi_j) + N(0, \sigma^2)$
  - The best choice of $\sigma$ is not 0!



$n = 1000, k = 100$

No noise: overestimate score by ≈0.10

Some noise: overestimate score by ≈0.06

average reported corr

noise scale (sigma) =

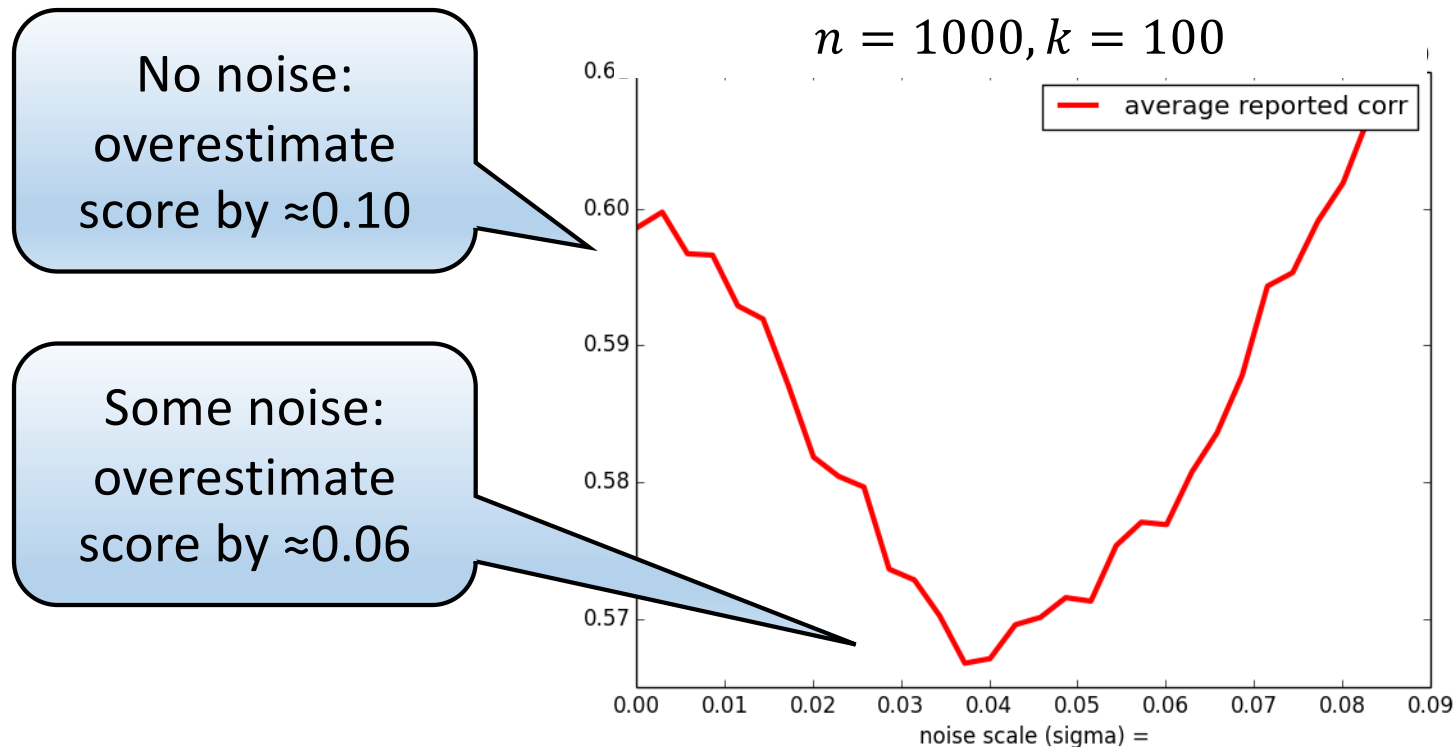# Case Study: ML Competitions

- **Improved estimator:** Add Gaussian noise $N(0, \sigma^2)$ to the estimated score of each classifier

  - Give answers $a_j = \text{score}_X(\varphi_j) + N(0, \sigma^2)$

  - The best choice of $\sigma$ is not 0!

> Minimized by
> $\sigma = \quad\quad$ ,
> achieving value

**Theorem** [DFHPRR'15, BNSSS**U**'16]**:** for appropriate $\sigma > 0$,
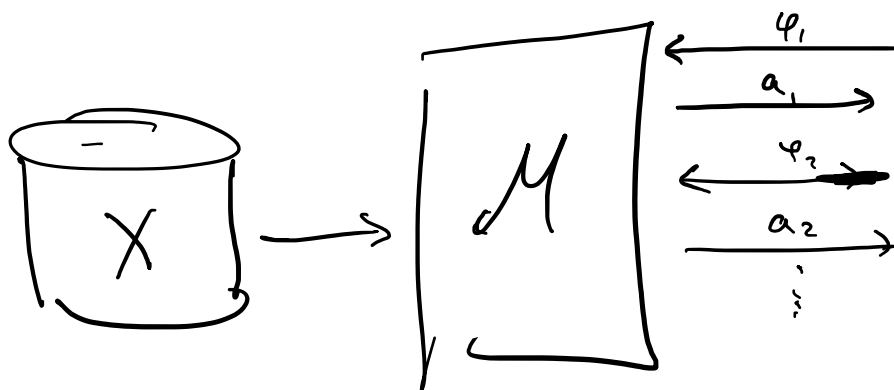
$$\mathbb{E}\left[\max_j a_j - \text{score}_P(\varphi_j)\right] \lesssim \frac{\sqrt{k}}{n\sigma} + \sigma$$

overfitting    noise

- Compare to $O\left(\sqrt{k/n}\right)$ when $\sigma = 0$

# General Setting



Queries: $\varphi : U \longrightarrow [0, 1]$

(data universe)

Desired answer: $\underset{X' \sim P}{\mathbb{E}}\left( \varphi(X') \right)$

Goal: minimize $\left\{ \underset{j}{\max} \left| a_j - \underset{X^* \sim P}{\mathbb{E}}\left( \varphi(X') \right) \right| \right.$

# Proof Overview

**Key Claim:** If $M'$ is an $(\varepsilon, \delta)$-DP mechanism that maps $X$ to a classifier, then $\mathbb{E}_{X,M}\big[\text{score}_X\big(M'(X)\big)\big] - \mathbb{E}_{X,M}\big[\text{score}_P\big(M'(X)\big)\big] = O(\varepsilon + \delta)$

How will we use this?

# Proof Overview

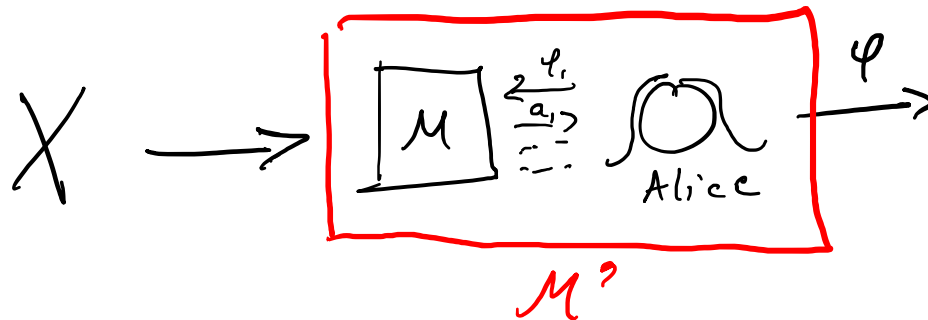**Key Claim:** If $M'$ is an $(\varepsilon, \delta)$-DP mechanism that maps $X$ to a classifier, then $\mathbb{E}_{X,M}\big[\text{score}_X(M'(X))\big] - \mathbb{E}_{X,M}\big[\text{score}_P(M'(X))\big] = O(\varepsilon + \delta) + \frac{1}{\sqrt{n}}$

How will we use this?



By POST-PROCESSING ($\forall$), $M'$ is $(\varepsilon, \delta)$-DP. if $M$ is $(\varepsilon, \delta)$-DP.

Say Alice is trying to find $\ell$ s.t. $\text{score}_X(\ell) \gg \text{score}_P(\ell)$

# Proof Overview

**Key Claim:** If $M'$ is an $(\varepsilon, \delta)$-DP mechanism that maps $X$ to a classifier, then $\mathbb{E}_{X,M}\big[\text{score}_X(M'(X))\big] - \mathbb{E}_{X,M}[\text{score}_P(M'(X))] = O(\varepsilon + \delta)$

- Proof Sketch:
  - Consider $(i, X_i, M'(X))$ and $(i, Z, M'(X))$ where $i \sim [n]$, $X \sim P^n, Z \sim P$ independently, and $M'$ is the mechanism

# Proof Overview

**Key Claim:** If $M'$ is an $(\varepsilon, \delta)$-DP mechanism that maps $X$ to a classifier, then $\mathbb{E}_{X,M}\left[\text{score}_X(M'(X))\right] - \mathbb{E}_{X,M}\left[\text{score}_P(M'(X))\right] = O(\varepsilon + \delta)$

- Proof Sketch:
    - Consider $\left(i, X_i, M'(X)\right)$ and $\left(i, Z, M'(X)\right)$ where $i \sim [n]$, $X \sim P^n, Z \sim P$ independently, and $M$ is the mechanism
        - **Sub-claim:** $\left(i, X_i, M'(X)\right) \approx_{\varepsilon, \delta} \left(i, Z, M'(X)\right)$
    - Observe that
        - $\mathbb{E}_{X,M}\left[\text{score}_X(M'(X))\right] = \mathbb{E}\left(f\left(i, X_i, M'(X)\right)\right)$
        - $\mathbb{E}_{X,M}\left[\text{score}_P(M'(X))\right] = \mathbb{E}\left(f\left(i, Z, M'(X)\right)\right)$
        - Where $f(i, y, m) = $ _____
    - **Fact:** If $A, B \in [0,1]$ satisfy $A \approx_{\varepsilon, \delta} B$, then $\mathbb{E}(A) \leq e^{\varepsilon}\mathbb{E}(B) + \delta.$

# Transfer Theorem

**Theorem:** Let $M$ be an $(\varepsilon, \delta)$-DP mechanism for answering a sequence of $k$ queries that is accurate on the sample, i.e.,

$$\Pr\left(\max_j \left|a_j - \text{score}_X(\varphi_j)\right| \leq \alpha\right) \geq 1 - \beta.$$
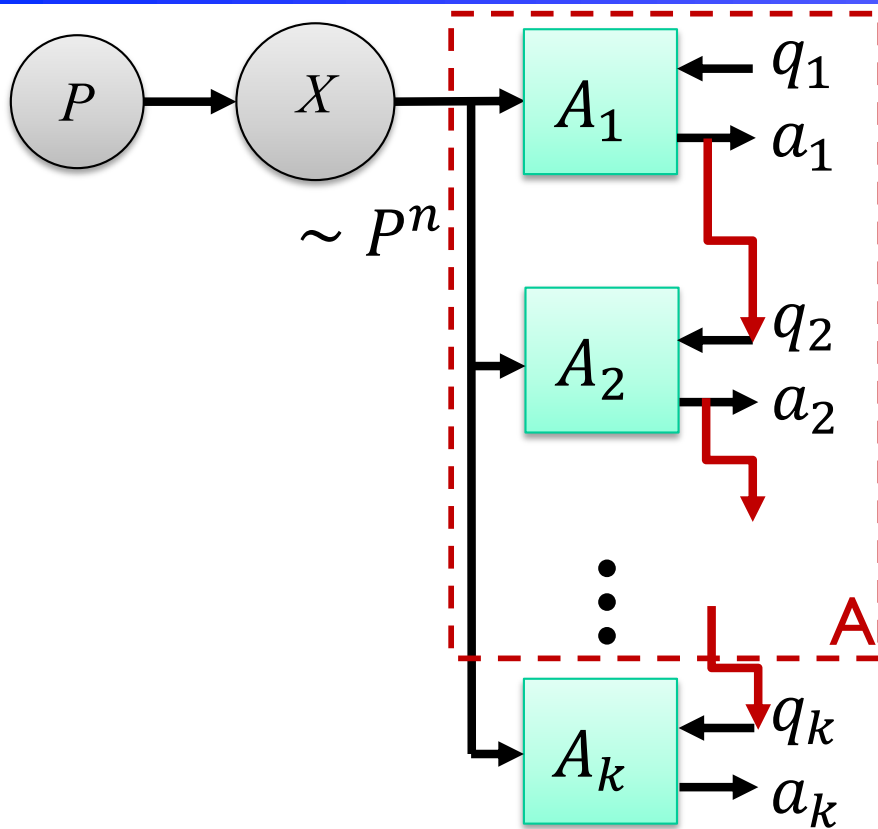
Then it is also accurate on the population:

$$\Pr\left(\max_j \left|a_j - \text{score}_P(\varphi_j)\right| \leq \alpha + \varepsilon + \sqrt{\beta} + \sqrt{\delta}\right) \gtrsim 1 - \sqrt{\beta} - \sqrt{\delta}.$$

This result is sufficient to analyze the Gaussian mechanism, as well as others based on MW-EM, for example

See Jung, Ligett, Neel, Roth, Sharifi-Malvajerdi, Moshe Shenfeld, *ITCS 2020* for a nice proof.
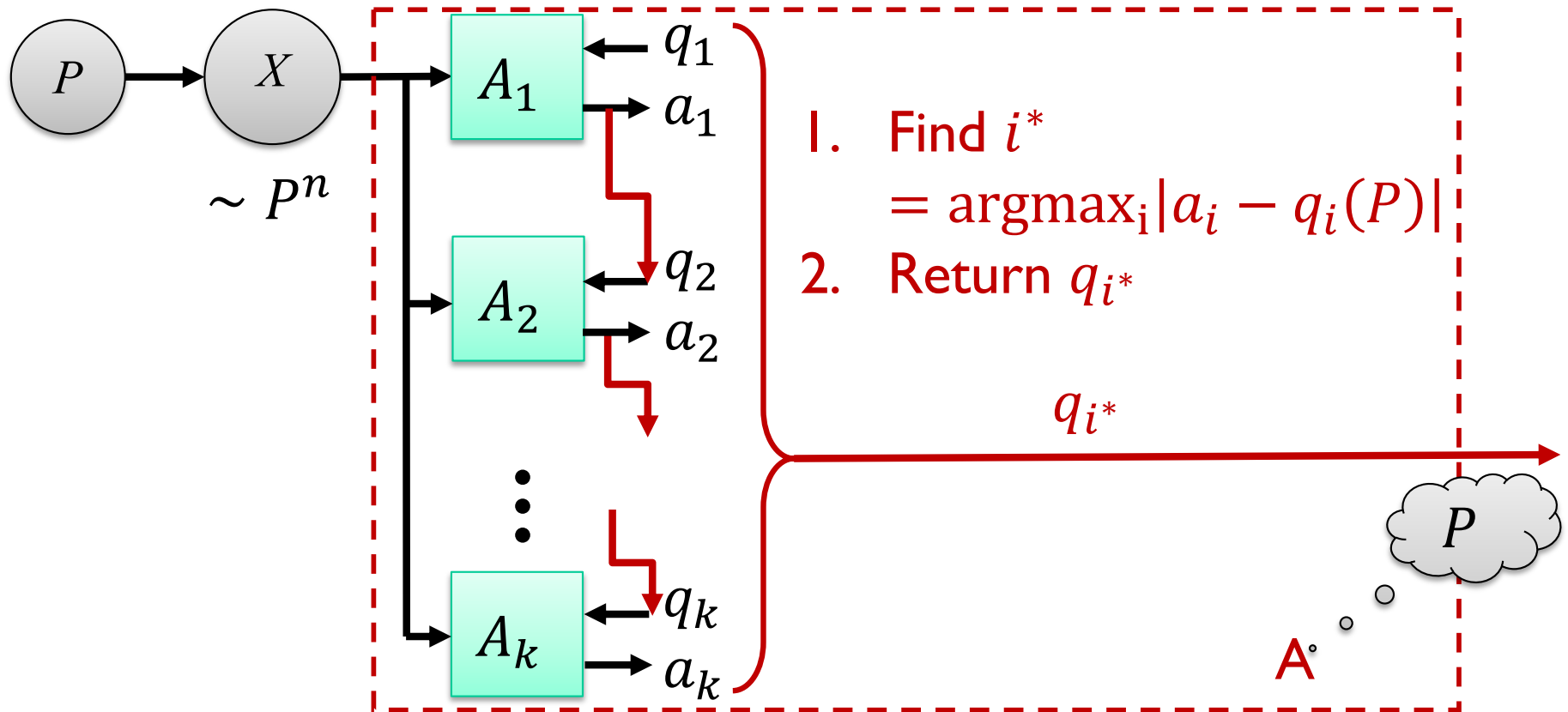
# What happens with Many Queries?

- Apply overfitting lemma at each round
  - ➤ Probability of overfitting adds up over rounds

# *"Monitor Argument"* [BNSSSU'16]



1. Find $i^*$
   $= \text{argmax}_i |a_i - q_i(P)|$
2. Return $q_{i^*}$

**Observation:**
$$\epsilon \geq Score_X(q_{i^*}) - Score_P(q_{i^*}) \geq \max_i |a_i - q_{i(P)}| - \alpha$$

- Stronger bounds
- Generalizes beyond linear queries

Versions based on ...

— other variants of DP.

— measures of information

$$I(X ; M(x))$$

(also other measures).

$\approx \varepsilon \sqrt{n} + (\cdots)$ when $M$ is DP and $X$ iid.

→ gives results for arbitrary hypothesis tests.