# CS541 Artificial Intelligence Guest Lecture on Mean Estimation
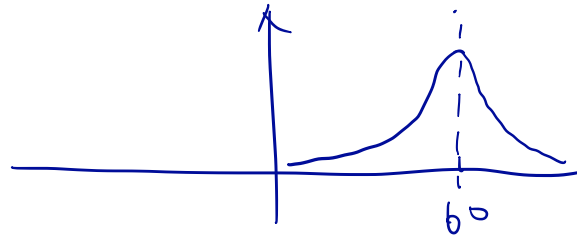
Lecturer: Shiwei Zeng
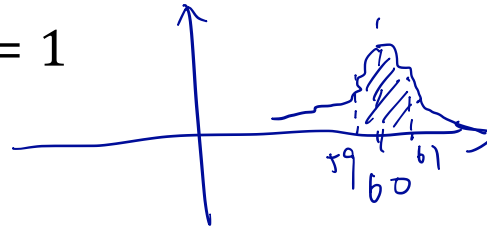
Szeng4 @ stevens. edu

# Estimating Average Height

- Assume $D = N(60, 1)$

- Assume $E[D] = 60, \text{Var}[D] = 1$

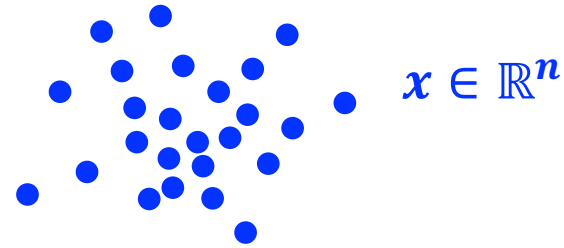- Estimator $\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i$ $\qquad S \sim D^m$

$$E[\hat{\mu}] = E_{S \sim D^m}\left[\frac{1}{n} \sum_{i=1}^{n} x_i\right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} E[x_i] = 60 = \text{ground truth}.$$

$$\text{Var}_{S \sim D^m}[\hat{\mu}] = \frac{1}{n^2} \sum_{i=1}^{n} \text{Var}[x_i] = \frac{1}{n^2} \cdot n = \frac{1}{n} \qquad n \uparrow$$

# ME in Higher Dimension

$$x \in \mathbb{R}^n$$

$$D$$

$$E[D] = ?$$

# When Data is Noisy

Adversary: Corrupt $\varepsilon$-fraction, $\varepsilon < \frac{1}{2}$

Total variation distance $D_1, D_2$

$$\frac{1}{2}\int |\phi_1 - \phi_2|\, dx = \frac{\varepsilon}{1-\varepsilon}$$

- 1-dimensional: (a lower bound)

$O(\varepsilon)$

$\phi_2$

$$\phi_1 \curvearrowleft D_1 = N(\mu_1, 1) \qquad D_2 = N(\mu_2, 1) \qquad |\mu_1 - \mu_2| \geq \Omega(\varepsilon)$$

$$Q_1, Q_2 \qquad D_\varepsilon = (1-\varepsilon)D_1 + \varepsilon \cdot Q_1 = (1-\varepsilon)D_2 + \varepsilon Q_2$$

$$Q_1 = \frac{1-\varepsilon}{\varepsilon}(\phi_2 - \phi_1) \cdot \mathbb{1}_{\phi_2 \geq \phi_1}$$

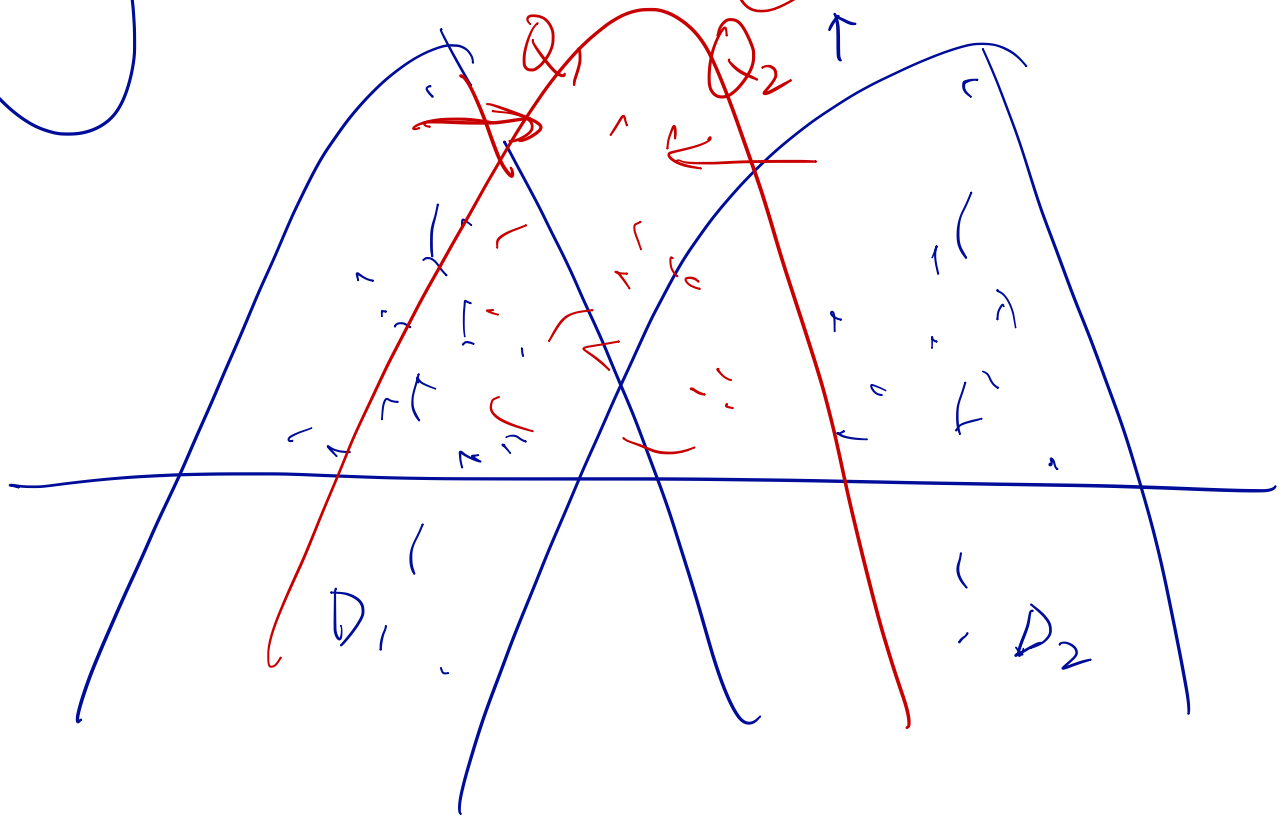$$Q_2 = \frac{1-\varepsilon}{\varepsilon}(\phi_1 - \phi_2) \cdot \mathbb{1}_{\phi_1 \geq \phi_2}$$

$$\text{Verify}: q_1 = (1-\varepsilon) \cdot \phi_1 + \left\{ \varepsilon \cdot \frac{1-\varepsilon}{\varepsilon}(\phi_2 - \phi_1) \cdot \mathbb{1}_{\phi_2 \geq \phi_1} \right. \qquad q_2 = (1-\varepsilon)\phi_2 + \varepsilon \frac{1-\varepsilon}{\varepsilon}(\phi_1 - \phi_2)$$

$$= \begin{cases} (1-\varepsilon) \cdot \phi_2 & \phi_2 \geq \phi_1 \\ (1-\varepsilon) \cdot \phi_1 & \phi_2 < \phi_1 \end{cases} \qquad = \begin{cases} \end{cases}$$

$$D_\epsilon = (1-\epsilon)D_1 + \boxed{\epsilon}\, Q_1$$

$$D_\epsilon = (1-\epsilon)D_2 + \boxed{\epsilon}\, Q_2$$

# Robust Mean Estimation

$[\frac{1}{2}, 1)$

**Mean Estimation**

$\mu_1$  $\mu_2$

$\varepsilon$-robust Mean Estimation

An $\varepsilon$ fraction is corrupted

$$D$$

$$D + D'$$

$\Omega(\varepsilon)$

bf 2016

$O(\varepsilon \cdot \sqrt{n})$

$1000 \cdot \varepsilon$

$n = 10^6$

$$E[D] = ?$$

$$E[D] = ?$$

$\Sigma = (0, \frac{1}{2})$

$\|\hat{\mu} - \mu\|_2 \geq 500$

# Natural approaches

- Learn each coordinate separately

$$|\hat{\mu} - \mu| \geq \Omega(\varepsilon)$$

$$\text{in } n\text{-dimension}: \quad \|\hat{\mu} - \mu\|_2^2 = \sum_{i=1}^{n} |\hat{\mu}_i - \mu_i|^2 \geq n \cdot \Omega(\varepsilon)^2 \geq \Omega(n\varepsilon^2)$$

$$\underline{\Omega(\sqrt{n} \cdot \varepsilon)}$$

$$e^n$$

# Natural approaches

- Maximum Likelihood Estimator

Negative Log likelihood = NLL

$$\min NLL(F, x_1, \cdots, x_m) = -\sum_{i=1}^{m} \log F(x_i)$$

$$F(x_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{\|x_i - \mu\|_2^2}{2}} \qquad \leftarrow \quad Var = 1$$

$$\min -\sum_{i=1}^{m} \log\left(\frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{\|x_i - \mu\|_2^2}{2}}\right)$$

$$\implies \min_{\mu \in \mathbb{R}^n} -\sum_{i=1}^{m} \left(\log\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{\|x_i - \mu\|_2^2}{2}\right)$$

$$\implies \arg\min \frac{1}{2}\sum_{i=1}^{m} \|x_i - \mu\|_2^2 = \hat{\mu} \longleftarrow \text{empirical mean.}$$

Can be quite bad.

Alg is not robust.

# Efficient Algorithm – Convex Programming

weight vector $\boxed{\hat{w}} = (w_1, w_2, \cdots, w_m)$

Goal: output $\underset{=}{\hat{w}}$, $\sum_{i=1}^{m} \hat{w}_i \cdot x_i = \hat{\mu} \longrightarrow \mu$.

min empirical variance.

s.t. $\boxed{\hat{w} \in W} \longleftarrow$ $O(n^6)$

# Efficient Robust Mean Estimation - Filter

$p(x) = x \cdot v^*$

$O(\sqrt{n})$

$D = N(\mu, I_n)$

1. Compute empirical mean and covariance $\mu_T, \Sigma_T$    $T$ : corrupted data set.

2. Compute largest eigenvalue $\lambda^*$ of $\Sigma_T - I$, and eigenvector $v^*$

3. If $\lambda^*$ is small, return $\mu_T$

   $\lambda^* \Sigma = v^* -$

4. Otherwise, find $t > C_1$ such that

$$\mathrm{Pr}_{X \in T}[|v^* \cdot (X - \mu_T)| > t] > C_2 e^{-t^2/2} + \frac{C_3 \varepsilon}{t^2 \log(n \log \frac{n}{\varepsilon \tau})}$$

5. Remove $X$ such that $|v^* \cdot (X - \mu_T)| > t$, go back to step 1

$O(\varepsilon)$

$\lambda^*$ : eigenvalue $\longrightarrow$ variance.

$$\lambda^* v^* = v^* \Sigma_T$$

$$\boxed{\begin{array}{l} \text{Var}\left[ x \cdot v^* \right] \\ x \sim N(0, I) \end{array}} = E\left[ (x \cdot v^*)^2 \right]$$

$$= \underset{x \sim D}{E}\left[ (v^{*T} x)(x^T \cdot v^*) \right]$$

$$= v^* \cdot \underbrace{E\left[ x x^T \right]}_{= \Sigma_T} v^*$$

$$= v^* \cdot \Sigma_T \cdot v^*$$

$$= \lambda^* v^* \cdot \underset{\longrightarrow = 1}{v^*}$$

$$\boxed{= \lambda^*}$$

# List-decodable Mean Estimation

$$\varepsilon \geq \frac{1}{2} \qquad \alpha = 1 - \varepsilon$$

**Mean Estimation**

<span style="color:red">List-Decodable</span> Mean Estimation

$$10^6 \not> 20\%$$
$$= 0.2 \times 10^6$$
$$200000$$

$$T$$

$$\rightarrow 20\% = \alpha$$

$$D$$

$$\mathbf{E[D] = ?}$$

$$\mathbf{E[D] \in \{\mu_1, \ldots, \mu_m\}}$$

# Algorithm: Multi-filtering

Gaussian Annulus Theorem.

$\sqrt{n}$

root = T

$T_i$

$\hat{\mu}_i$

$T_i \cap T_2 = \emptyset$

$T_i \, T_2$

$\left(\frac{1}{\alpha}\right)$

- A tree of subsets $T_i$'s,

① Clustering

$x_6$

$x_i^2 \; x_1 x_2 x_3$

$\mathbb{E}[XX^T]$

$[x \; x_1 \; x_1 x_2 \; x_3 \ldots]$

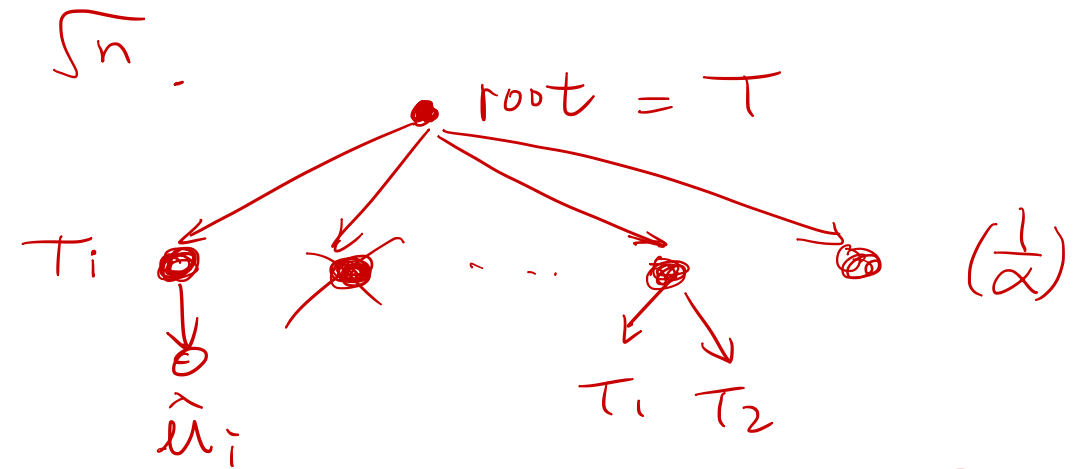② • Iterate through each node
  - (1) Create a leaf node, an estimate $\hat{\mu}_i$
  - (2) Create child nodes, subsets $T_i$'s
    - a.  One node, cleaner set
    - b.  Two nodes, overlapping subsets
  - (3) Delete if it can't be $\alpha$-good.

$\alpha$-fraction

- No more filtering, then return all $\hat{\mu}_i$'s

$O\left(\sqrt{\log \frac{1}{\alpha}}\right)$

$O\left(\frac{1}{\alpha^{\frac{1}{2d}}}\right)$  $d \in \mathbb{N}^+$, d: degree of polynomial

$T_i: \; \alpha T_i \longrightarrow$ good samples.

$err = \min_i \|\hat{\mu}_i - \mu\|_2 \leq O\left(\frac{1}{\sqrt{\alpha}}\right)$

11