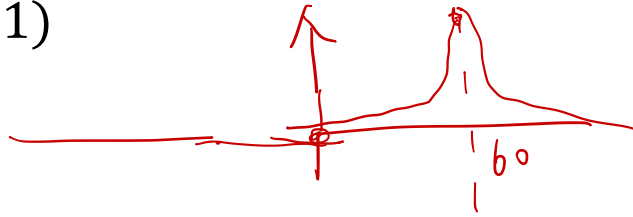# CS541 Artificial Intelligence Guest Lecture on Mean Estimation
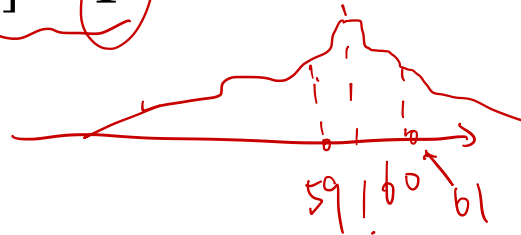
Lecturer: Shiwei Zeng

# Estimating Average Height

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

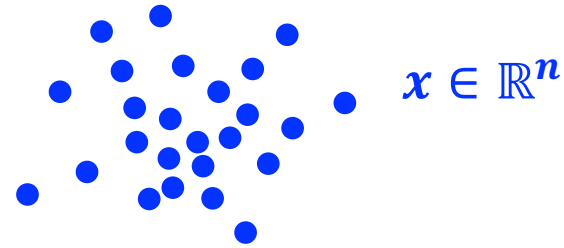- Assume $E[D] = 60, \text{Var}[D] = 1$

- Estimator $\hat{\mu} = \frac{1}{n}\sum_{i=1}^{n} x_i$

$x_i \sim D$ i.i.d.

$$E[\hat{\mu}] = E\left[\frac{1}{n}\sum_{i=1}^{n} x_i\right] = \frac{1}{n}\sum_{i=1}^{n} E[x_i] = \mu$$

$$\text{Var}[\hat{\mu}] = \text{Var}\left[\frac{1}{n}\sum_{i=1}^{n} x_i\right] = \frac{1}{n^2}\sum_{i=1}^{n} \text{Var}[x_i] = \frac{1}{n}$$

# ME in Higher Dimension



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

$$D$$

$$E[D] = ?$$

# When Data is Noisy

$O(\cdot)$

$\varepsilon < \frac{1}{2}$   $\varepsilon = 3.0\%$

$\varepsilon \cdot c$   $c > 0$

- 1-dimensional: (a lower bound)

$\|\hat{\mu} - \mu\|_2^2 \overset{``=''}{\geqslant} \Omega(\varepsilon)$ if corrupt $\varepsilon$-fraction.

$D_1 = N(\mu_1, 1)$   $D_2 = N(\mu_2, 1)$   Construct $Q_1$ and $Q_2$ s.t.

$\|\mu_1 - \mu_2\| = \Omega(\varepsilon)$ and

$\boxed{D_\varepsilon} = (1-\varepsilon) D_1 + \varepsilon Q_1 = (1-\varepsilon) D_2 + \varepsilon Q_2$ (1)

Let $\phi_1$ be pdf of $D$, $\phi_2$: pdf of $D_2$. Let $\mu_1, \mu_2$ be s.t. the total variance distance

between $D_1, D_2$ is

$$\frac{1}{2}\int |\phi_1 - \phi_2| dx = \frac{\varepsilon}{1-\varepsilon} \implies \|\mu_1 - \mu_2\| \geqslant \frac{2\varepsilon}{1-\varepsilon}.$$

$Q_1 = \frac{1-\varepsilon}{\varepsilon}(\phi_2 - \phi_1) \cdot 1_{\phi_2 > \phi_1}$ and $Q_2 = \frac{1-\varepsilon}{\varepsilon}(\phi_1 - \phi_2) 1_{\phi_1 > \phi_2}$.

4

$$D_1 = N(\mu_1, 1) \qquad D_2 \; N(\mu_2, 1)$$

$$\|\mu_1 - \mu_2\| = \Omega(\varepsilon) \quad \text{and}$$

$$\boxed{D_\varepsilon} = (1-\varepsilon)D_1 + \varepsilon Q_1 = (1-\varepsilon)D_2 + \varepsilon Q_2 \qquad (1)$$

Let $\phi_1$ be pdf of $D$, $\phi_2$ pdf of $D_2$. Let $\underline{\mu_1, \mu_2}$ be s.t. the total variance distance between $D_1, D_2$ is

$$\frac{1}{2}\int |\phi_1 - \phi_2| \, dx = \frac{\varepsilon}{1-\varepsilon} \implies \|\mu_1 - \mu_2\| \geq \frac{2\varepsilon}{1-\varepsilon} \geq \Omega(\varepsilon)$$

$$4\varepsilon \geq \qquad \geq 2\varepsilon \qquad \qquad \in (\tfrac{1}{2}, 1)$$

$$Q_1 = \frac{1-\varepsilon}{\varepsilon}(\phi_2 - \phi_1) \cdot \mathbb{1}_{\phi_2 \geq \phi_1} \quad \text{and} \quad Q_2 = \frac{1-\varepsilon}{\varepsilon}(\phi_1 - \phi_2)\mathbb{1}_{\phi_1 \geq \phi_2}.$$

$$(1-\varepsilon)\phi_1 + \varepsilon \cdot \frac{1-\varepsilon}{\varepsilon}(\phi_2 - \phi_1)\cdot \mathbb{1}_{\phi_2 \geq \phi_1} \qquad \mathbb{1}_{\phi_2 \geq \phi_1} = \begin{cases} 1 & \phi_2 \geq \phi_1 \\ 0 & \phi_2 < \phi_1 \end{cases}$$

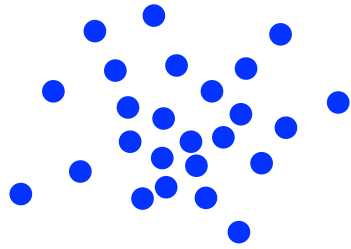$$= (1-\varepsilon)\phi_1 + (1-\varepsilon)(\phi_2 - \phi_1) \cdot \mathbb{1}_{\phi_2 \geq \phi_1}$$

$$= \begin{cases} (1-\varepsilon)\cdot \phi_2 & \phi_1 \leq \phi_2 \\ (1-\varepsilon)\cdot \phi_1 & \phi_1 > \phi_2 \end{cases}$$

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

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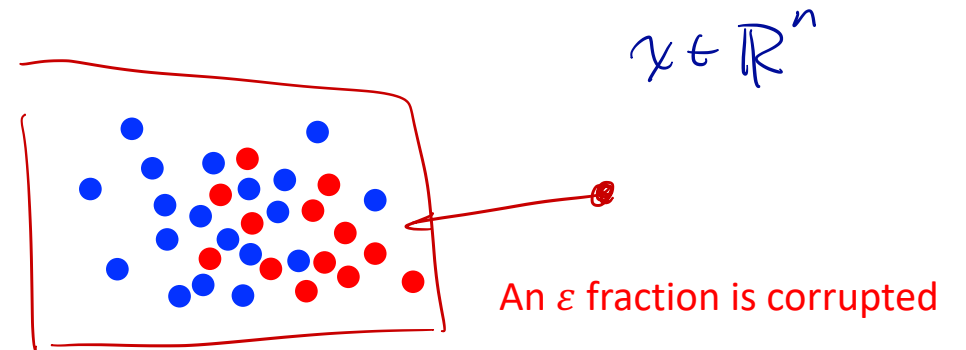
# Robust Mean Estimation

**Mean Estimation**

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

**$\varepsilon$-robust** Mean Estimation

$D$

$D + D'$

An $\varepsilon$ fraction is corrupted

$E[D] = ?$

$E[D] = ?$

# Natural approaches

- Learn each coordinate separately

$$\hat{\mu} \quad s.t. \quad \|\hat{\mu}_i - \mu_i\| \geq \Omega(\varepsilon).$$

in $n$-dim.

$$\|\hat{\mu} - \mu\|_2^2 \geq \sum_{i=1}^{n} (\hat{\mu}_i - \mu_i)^2 = n \cdot \varepsilon^2$$

$$\|\hat{\mu} - \mu\|_2 \geq \sqrt{n} \cdot \varepsilon$$

$$n = 10^6$$

$$\sqrt{n} \cdot \varepsilon = 10^3 \cdot \varepsilon \qquad \varepsilon (0, \tfrac{1}{2})$$

# Natural approaches

$e^n$

- Maximum Likelihood Estimator

Negative Log Likelihood

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

$$S = \{x_i\}_{i=1}^{m}$$

Assume $F$ is Gaussian: $\quad F(x_i) = \dfrac{1}{\sqrt{2\pi}} \cdot e^{-\frac{\|x_i - \mu\|_2^2}{2}}$

$$\min \; NLL = \min - \sum_{i=1}^{m} \log \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{\|x_i - \mu\|_2^2}{2}}$$

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

$$\implies \min_{\mu} \frac{1}{2} \sum_{i=1}^{m} \|x_i - \mu\|_2^2$$

$\hat{\mu}$ : empirical mean

$$= \frac{1}{m} \sum_{i=1}^{m} x_i$$

7

# Efficient Algorithm – Convex Programming

$\sim 7:30$

Output a $\hat{w} = (w_1, w_2, \cdots w_m)$ such that

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^{m} w_i x_i \quad \text{is close to } \mu.$$

by solving a convex program.

$$O(n^6)$$

$D = N(\underline{\mu}, I)$  $\forall v \in \mathbb{R}^n$  $\Pr_{x \sim D}[|v \cdot (x - \mu)| \geq t] \leq \exp\left(-\frac{t^2}{2}\right)$  ☆  $\varepsilon(0, \frac{1}{2})$  $\underline{\tau = 0.01}$  $0.99$

# Efficient Robust Mean Estimation - Filter

$\tau$: corrupted dataset.

$\Pr_{x \sim D}[|p(x) - E[p(x)]| \geq t] \in \boxed{\phantom{x}}$  $\frac{1}{m}\sum_{i=1}^{m} x_i$  $\frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_T)(x_i - \mu_T)^T$

1. Compute empirical mean and covariance $\mu_T$, $\Sigma_T = I$

2. Compute largest eigenvalue $\lambda^*$ of $\boxed{\Sigma_T - I,}$ and eigenvector $v^*$
   $= 0$

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

4. Otherwise, find $t > C_1$ such that  $C_1, C_2, C_3 > 0$  $\tau$  $= \delta$ dependent on $m$ ↑

$$\boxed{\Pr_{X \in T}[|v^* \cdot (X - \mu_T)| > t]} > C_2 e^{-t^2/2} + \boxed{\frac{C_3 \varepsilon}{t^2 \log(n \log \frac{n}{\varepsilon \tau})}}$$
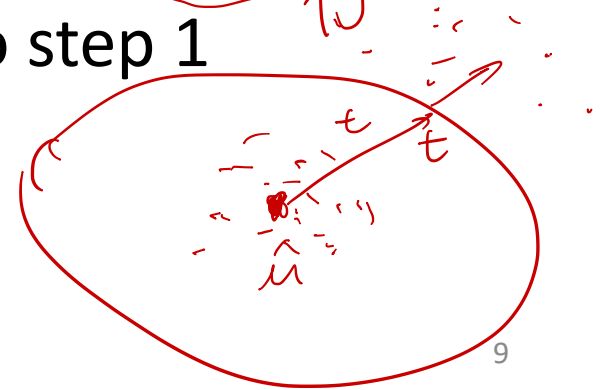
5. Remove $X$ such that $|v^* \cdot (X - \mu_T)| > t$, go back to step 1
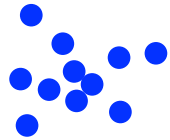
$\sqrt{n}$, n-dim

$\hat{\mu} \to \mu$

$\|\hat{\mu} - \mu\|_2 \leq O(\varepsilon) \leq 1$

# List-decodable Mean Estimation

$T$: corrupted data set.

**Mean Estimation**
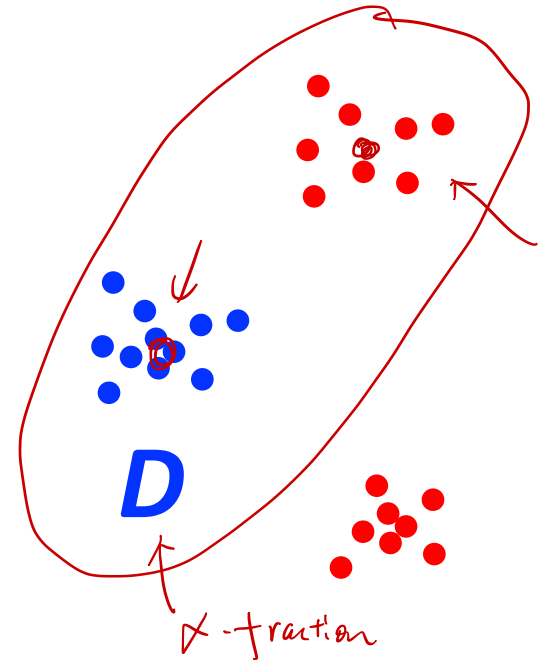
List-Decodable Mean Estimation

$\alpha$: fraction of clean samples.

$D$

$E[D] = ?$

$k$ golden queries
→ a list of $\ell^k$ estimated meas.

$x_1$

$D$

$\alpha$-fraction

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

$\log(m)$ queries

$m = O\left(\frac{1}{\alpha}\right)$

# Algorithm: Multi-filtering

*Handwritten (top):* $T_i$ is $\alpha$-good = $\alpha$-fraction of $T_i$ are clean.

- A tree of subsets $T_i$'s,

- 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.

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

*Handwritten (right):*

$T$ → $T_1, T_2, T_3, \ldots, T_m$    ← $O\left(\frac{1}{\alpha}\right)$

list ← $\{\mu_1 \ldots \mu_m\}$

$T \to T_1, T_2$

$T_1 \cap T_2 = $ ball with radius $\sqrt{n}$.

$80\%$ $T_1$    $T_2$
$50\%$    $50\%$

$\alpha$-good $T_i$ = $\alpha$-fraction of $T_i$.

$\hat{\mu}$    $\mu$    $O\left(\frac{1}{\alpha^{\frac{1}{2}}}\right)$

$poly\left(\frac{1}{\alpha}\right) \Rightarrow O\left(\frac{1}{\alpha^{\frac{1}{2d}}}\right)$    $d$: degree

*Handwritten (bottom left):*

$\min_{i \in S\{\mu_1 \ldots \mu_m\}} \| \hat{\mu}_i - \mu \|_2 = O\left(\frac{1}{\sqrt{\alpha}}\right) = 1000$

$\alpha = 10^{-6}$