## Lecture #10: Measurement Error

- R-55: measurement error

$$\begin{bmatrix} 16\ oz \\ 16 \\ \vdots \\ 16 \end{bmatrix} \xrightarrow[\text{for more sig figs}]{\text{turn dial}} \begin{bmatrix} 16.0\ oz \\ 16.0 \\ \vdots \\ 16.0 \end{bmatrix} \rightarrow \begin{matrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{matrix} \begin{bmatrix} 15.97 \\ 16.01 \\ \vdots \end{bmatrix} \Big\} n$$

deterministic

probabilistic (stochastic)

- Basic measurement error model

$y_1 = (\text{true value}) + (\text{bias}) + (\text{random error})_1$

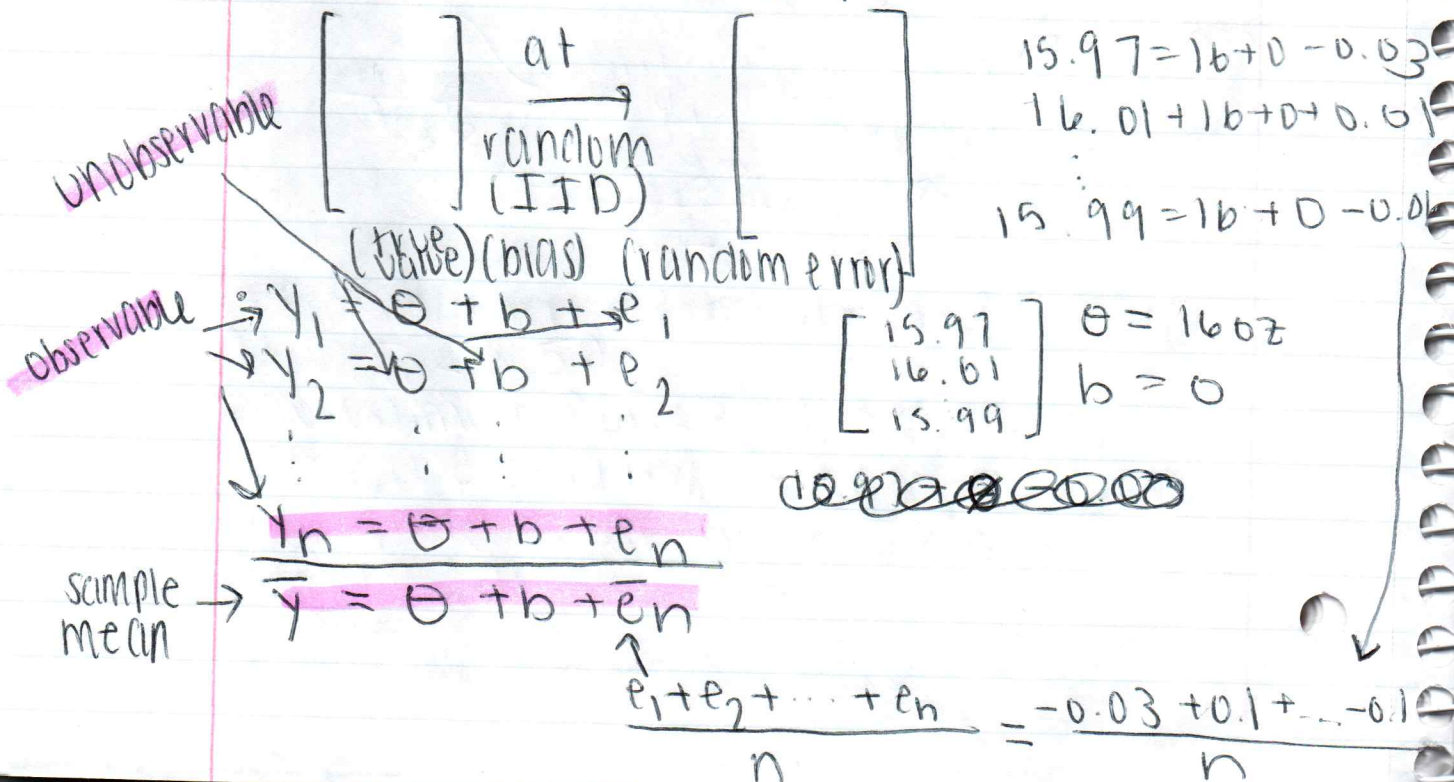$y_2 = (\text{true value}) + (\text{bias}) + (\text{random error})_2$

$\vdots$

$y_n = (\text{true value}) + (\text{bias}) + (\text{random error})_n$

IID mean $=0$ SD $=\sigma$

- Bias is a systematic tendancy to over or underestimate the truth

- Unbiased $= (\text{bias}=0)$ (no bias)

pop                      sample

$$\begin{bmatrix} \ \\ \ \\ \ \\ \ \end{bmatrix} \xrightarrow[\substack{\text{random} \\ \text{(IID)}}]{\text{at}} \begin{bmatrix} \ \\ \ \\ \ \\ \ \end{bmatrix}$$

unobservable

$15.97 = 16 + 0 - 0.03$

$16.01 + 16 + 0 + 0.01$

$\vdots$

$15.99 = 16 + 0 - 0.01$

(value)(bias) (random error)

observable $\rightarrow y_1 = \theta + b + e_1$

$y_2 = \theta + b + e_2$

$\vdots \quad \vdots \quad \vdots$

$\begin{bmatrix} 15.97 \\ 16.01 \\ 15.99 \end{bmatrix}$  $\theta = 16\ oz$  $b = 0$

$y_n = \theta + b + e_n$

sample mean $\rightarrow \bar{y} = \theta + b + \bar{e}n$

$$\frac{e_1 + e_2 + \cdots + e_n}{n} = \frac{-0.03 + 0.1 + \cdots - 0.1}{n}$$

- cancellation of ⊕ and ⊖ errors will yield an $\bar{e}_n$ that is highly likely to be closer to 0 than any of the errors $e_1, ... e_n$ themselves

$$\bar{y}_n = \theta + b + \bar{e}_n$$

(sample mean) = (truth) + (bias) + (mean random errors)

- as $n\uparrow$ $\bar{e}_n \to 0$ highly likely
- Therefore as $n\uparrow$, $\bar{y}_n \to \theta + b$
- $\bar{y}_n \to$ (truth) $\theta$, only when bias = 0
  - "get more good data"
  
  unbiased ↗

- 1936 FDR vs. Landon

| Literary digest | 12 mil letters |
| --- | --- |
| | 2.5 mil replies |

result: Landon 60%, FDR 40% · prediction
actual: FDR 60%, Landon 40%
-20% point error

1 = FDR   0 = Landon
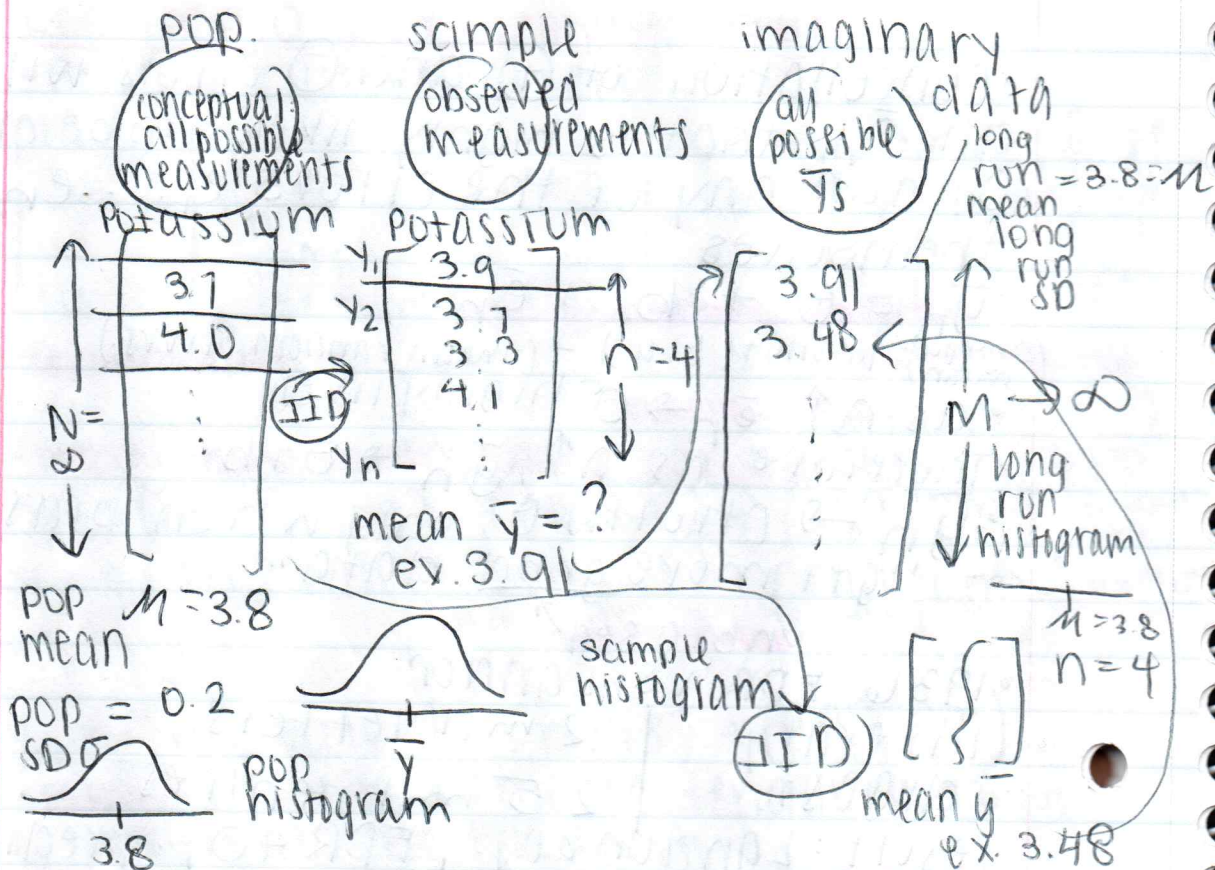
$$\bar{y}_n = \theta + b + \cancel{\bar{e}_n}$$

$\underbrace{\qquad\qquad}$
60% (-20%)

$$\begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \begin{bmatrix} 1s \\ 0s \end{bmatrix} \quad n = 2.5\,mil$$
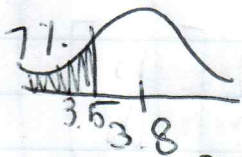
mean 40%

pop.
(n | θ) →want→ (sample θ)
representative

- Phonebooks, club membership (country clubs)
  - oversampled rich ppl → Republican

# probability model for measurement error

| POP. | sample | imaginary |
|---|---|---|
| (conceptual) all possible measurements | observed measurements | (all possible) $\bar{Y}$s |

imaginary data

long run = $3.8 = \mu$ mean long run SD

potassium

$$\begin{bmatrix} 3.7 \\ 4.0 \\ \vdots \end{bmatrix}$$

$N = \infty$ ↓

potassium

$$\begin{bmatrix} y_1 & 3.9 \\ y_2 & 3.7 \\ & 3.3 \\ & 4.1 \\ \vdots \\ y_n \end{bmatrix}$$  (IID)  $n = 4$

$$\begin{bmatrix} 3.91 \\ 3.48 \\ \vdots \\ \vdots \end{bmatrix}$$

mean $\bar{y} = ?$
ex. 3.91

$M \to \infty$

long run histogram

$\mu = 3.8$
$n = 4$

POP mean $\mu = 3.8$

POP $= 0.2$ SD

pop histogram

$\bar{y}$

sample histogram

(IID)  $\left[\begin{array}{c} \xi \end{array}\right]$  mean $\bar{y}$ ex. 3.48

$P(\text{misclassification w/ } n=1) = 7\%$

↑ too high

$SD = 0.2$

7%

hist of $y_1$

3.5 3.8

$\sim \dfrac{3.5 - 3.8}{0.2} = -1.5$

P(misclassification w/ n=4)?
$= P(\bar{y} < 3.5)$

• To estimate $P(\bar{y} < 3.5)$ we have to imagine getting lots of $\bar{y}$s and compute % of time $\bar{y} < 3.5$ in those repetitions

• Expected value $\bar{Y} = EV$ of $\bar{Y} = E_{IID}(\bar{y}) = \mu = 3.8$