

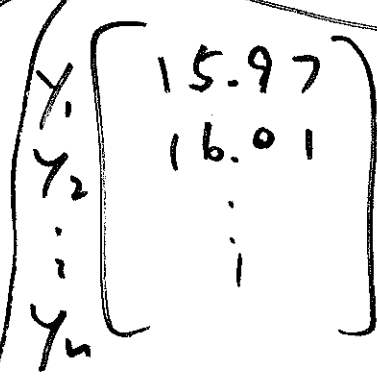
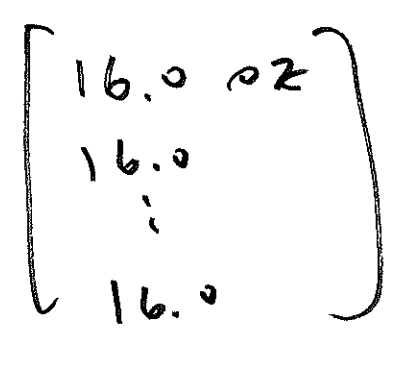
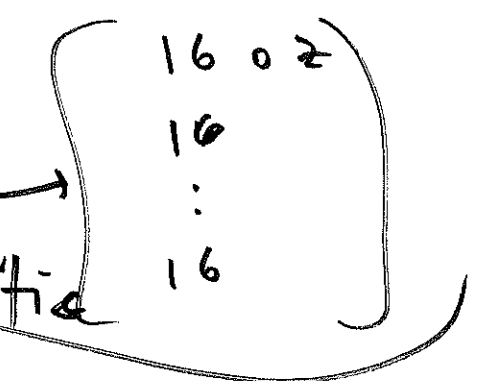
this measurement error
 time:
 next time: statistical inference

read: JD(B) ch. 11
 LN pp. 137-156
 AMS 30 Oct 18

today: LN pp. 127 → 1

R (55)

measurement error



basic measurement error model

$$\begin{aligned}
 y_1 &= (\text{true value}) + (\text{bias}) + \left(\begin{array}{c} \text{IID mean 0} \\ \text{SD } \sigma \\ \text{random} \\ \text{error} \end{array} \right)_1 \\
 y_2 &= (\text{true value}) + (\text{bias}) + \left(\begin{array}{c} \text{random} \\ \text{error} \end{array} \right)_2 \\
 &\vdots \\
 y_n &= (\text{true value}) + (\text{bias}) + \left(\begin{array}{c} \text{random} \\ \text{error} \end{array} \right)_n
 \end{aligned}$$

bias a systematic tendency to over- or underestimate the truth

unbiased = (bias = 0) (no bias)



observable

unobservable

(n)	(true)	(bias)	(random errors)
y_1	θ	b	e_1
y_2	θ	b	e_2
\vdots	\vdots	\vdots	\vdots
y_n	θ	b	e_n

(3)

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

sample mean \nearrow

\nearrow

$$\frac{e_1 + e_2 + \dots + e_n}{n}$$

$\begin{bmatrix} 15.97 \\ 16.01 \\ \vdots \\ 15.99 \end{bmatrix}$	$\theta = 16.02$ $b = 0$	$15.97 = 16 + 0 + (-0.03)$ $16.01 = 16 + 0 + (0.01)$ \vdots $15.99 = 16 + 0 + (-0.01)$
---	-----------------------------	---

$$\frac{e_1 + e_2 + \dots + e_n}{n} = \frac{(-0.03) + (0.01) + \dots + (-0.01)}{n}$$

cancellation of \oplus, \ominus errors will yield an \bar{e}_n that is ^{highly likely to} be closer to 0 than any of the errors e_1, \dots, e_n themselves

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

(sample mean) = (truth) + (bias) + (mean of window errors)

as $n \uparrow$ \bar{e}_n ~~is highly likely~~ ^{is highly likely} to $\rightarrow 0$

therefore as $n \uparrow$, $\bar{y}_n \rightarrow \theta + b$

$\bar{y}_n \rightarrow$ (truth) θ only when

bias = 0

"get more good data"

1936

FJR vs. Lardner

unbiased

Literary Digest

12 million letters

2.5 million replies

result: London 60% FDR 49% 5
 prediction

actual FDR 60% London 40%

20 percentage point error

1 = FDR
 0 = London

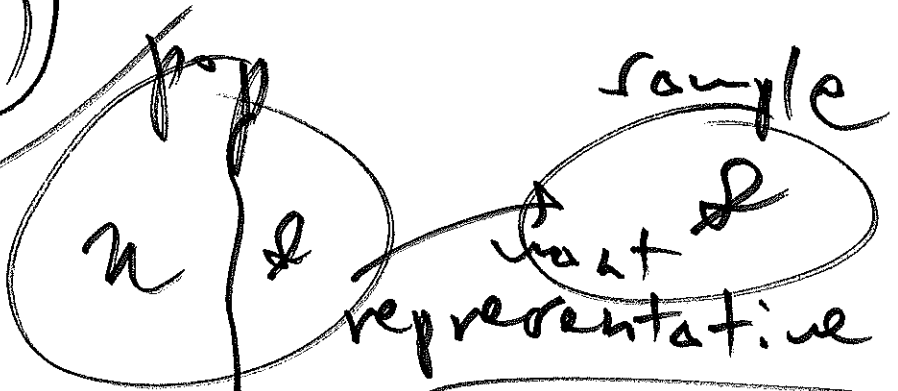
$y_i = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ $h = 2.5$ million
 near 40%

$\hat{y}_n = \theta + b + \tilde{e}_n$

$40\% = 60\% + (-20\%)$

phone books

club membership lists
 country clubs



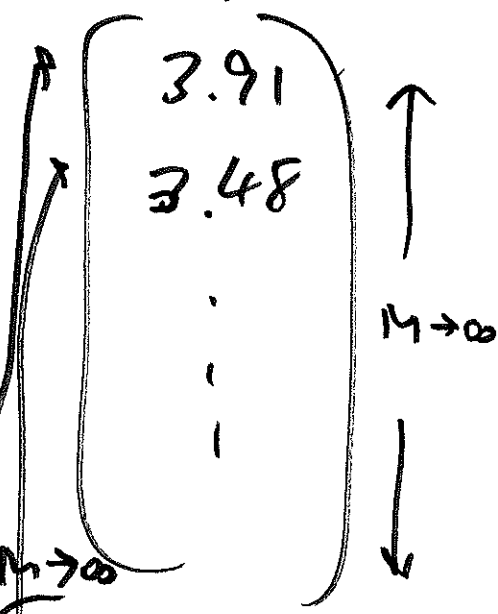
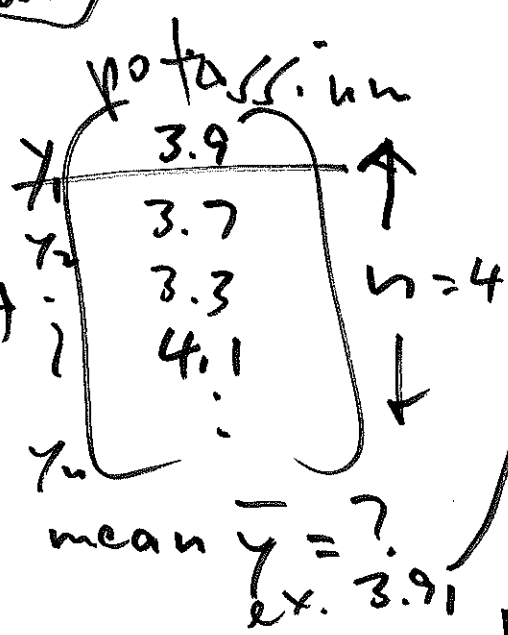
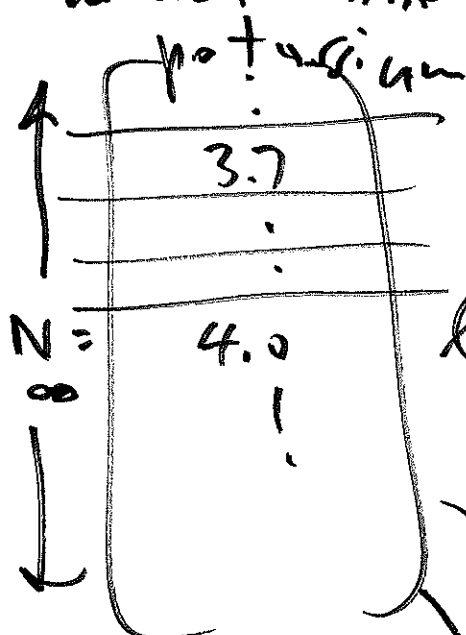
oversampled
 rich people + R

pop
conceptual:
all possible
measurements

prob.
model
for
meas.
error

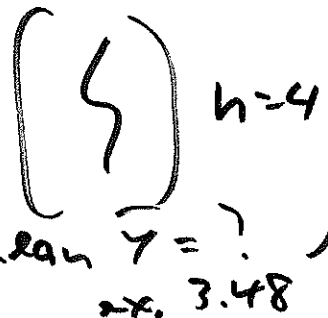
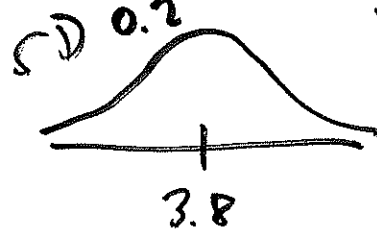
sample
the observed
measurements

imag.
data
all possible
 \bar{y}_s

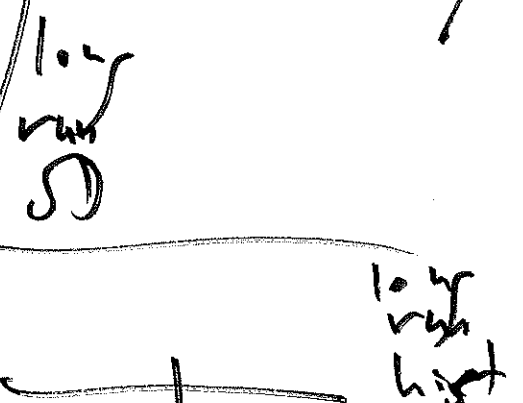


pop mean $\mu = 3.8$

pop SD $\sigma = 0.2$

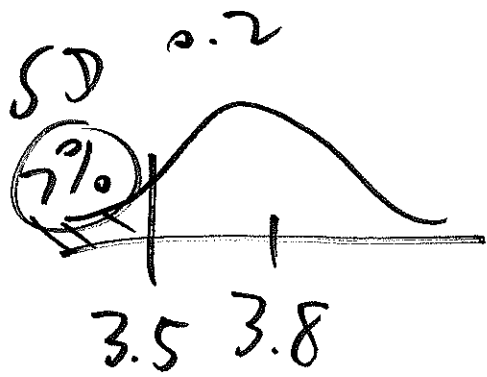


low var
mean $\bar{y} = 3.8 = \mu$



to estimate $P(\bar{y} < 3.5)$, we have to
imagine getting lots of \bar{y}_s &
compute % of time $\bar{y} < 3.5$ in those repetitions

$P(\text{misclassification with } n=1)$ ⑦



hist
of
 \bar{y}_1

(12.4%) (too high)

$P(\text{misclassification with } n=4)$

= ?

= $P(\bar{y} < 3.5)$

expected

value
of \bar{y} = EV of \bar{y}

$$= E_{\text{IID}}(\bar{y}) = \mu = 3.8$$