

STAT 131
27 May 20

read JS ch. 5

this useful
time: distributions
next time: catalog

people in U.S. = 325M

Catch-up lecture

eligible voters = 120M = T

& motivated

θ = population proportion favoring Clinton, 2 weeks before election

if can achieve random or like-at-random sampling from p

pop. all voters favor C?

the observed sample favor C?

$T = \begin{pmatrix} 1s \\ 2s \\ 0s \end{pmatrix}$ mean θ

like at random mean θ

$S = \begin{pmatrix} 1s \\ 2s \end{pmatrix}$ mean $\hat{\theta}$

IID

$$(S | n, \theta) \sim \text{Binomial}(n, \theta) \quad (2)$$

$$E(\hat{\theta}) = E\left(\frac{S}{n}\right) = \frac{1}{n} \underbrace{E(S)}_{n\theta} = \theta$$

So $\hat{\theta}$ is an unbiased estimate of θ

$$V(\hat{\theta}) = V\left(\frac{S}{n}\right) = \frac{V(S)}{n^2}$$

$$= \frac{n\theta(1-\theta)}{n^2} = \frac{\theta(1-\theta)}{n} \quad \text{and}$$

$$SD(\hat{\theta}) = \sqrt{\frac{\theta(1-\theta)}{n}} \quad \downarrow \text{or } \frac{1}{\sqrt{n}}$$

$$= O\left(\frac{1}{\sqrt{n}}\right)$$

uncertainty about θ based on $\hat{\theta}$

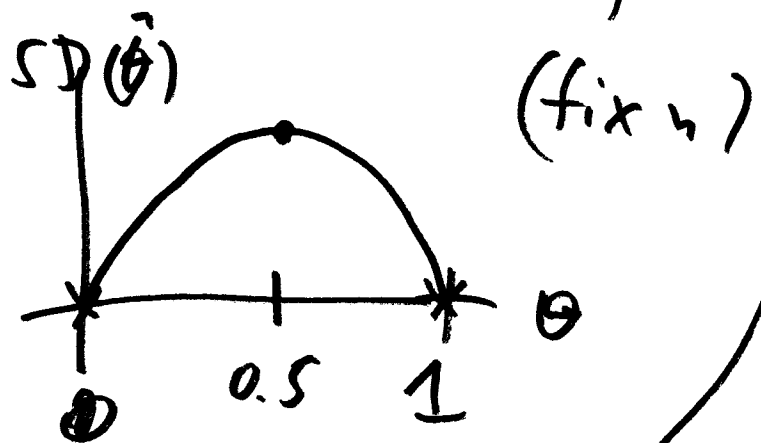
square root law

to cut $SD(\hat{\theta})$ in half, need
to quadruple the sample size

Q: How big does n have to
be to get $SD_{IID}(\hat{\theta}) = 1\%$?

A: $SD_{IID}(\hat{\theta}) = \sqrt{\frac{\theta(1-\theta)}{n}} = .01$

given 2 candidates & polarization
of electorate, $\theta = 0.5$



$SD_{IID \text{ worst case}}(\hat{\theta}) = \sqrt{\frac{0.5 \cdot (1-0.5)}{n}}$
 $= \frac{0.5}{\sqrt{n}}$

$$\frac{0.5}{\sqrt{n}} = .01 \rightarrow n = \left(\frac{0.5}{0.01}\right)^2 = 2500$$

T = 120M

$n_{\text{worst}} =$ n_{best}
 Case IID Case SRS

big issue,

however : really hard to

achieve like-at-random

~~(T, C)~~ (T, C, ^{no} response)

error due to random
 unlucky sampling

$$RMSE(\hat{\theta}) = \sqrt{V_{IID}(\hat{\theta}) + (\text{bias}(\hat{\theta}))^2}$$

if $\hat{\theta}$ is biased = $SD_{IID}(\hat{\theta})$ with no bias

Literary Digest (LD) poll

1936 $\left\{ \begin{array}{l} (D) \text{ FDR} \\ (R) \text{ Alf Landon} \end{array} \right.$

(5)

in summer of 1936, mailed out 10M
letters, only 2.27 M responded,
~~rest~~: $\hat{\theta} = 60\%$ Landon; truth:
estimate 59% R (!)

error: 19 percentage points (!)

how is this possible with $n = 2.27 \text{ M}$?

$$\text{SD}(\hat{\theta}) = \sqrt{\frac{(0.6)(0.4)}{2.27 \text{ M}}} = .000325$$

\leftarrow their mistake $\approx .003\%$

$$\text{RMSE}(\hat{\theta}) = \sqrt{\underbrace{V(\hat{\theta})}_{\text{IIS}} + \underbrace{(\text{bias}(\hat{\theta}))^2}_{.19}}$$

\leftarrow biased

George Gallup: correctly estimated ^⑥

LD result to within 1% percentage

point before election, correctly

estimated FDR's ^{vote} share was

about 59%

LD sources of bias

① non-response bias: people who

chose to mail postcard back

were systematically different in

voting behavior than people who

chose not to

② how obtain 10M ^{valid} addresses in

1936? phonebooks + club membership lists

1M subscribers

Country clubs, ...

not everybody had phone,
bias toward wealthy → London (R)

bias toward wealthy → London (R)

$$RMSE(\hat{\theta}) = \sqrt{V_{IID}(\hat{\theta}) + [\text{bias}(\hat{\theta})]^2}$$

→ 0 as n ↑

typical amount by which $\hat{\theta}$ misses θ

real enemy in (data) science

1
A sample
is representative
of the population
from which
it's taken

⑧
The sample
is
like-IID
or at least
like-at-random