

SRS: design based theory

Week 2 (2.8)

Stat 260, St. Clair

"Intro stats" vs. "design-based"

An Intro Stats story: Responses are randomly selected from a "normal population"

- the population is infinite
- the observed value Y_i is **random**: $Y_i \sim N(\mu, \sigma^2)$

"Intro stats" vs. "design-based"

Design-based story:

- Responses are fixed values (not random, just unknown)
- the population is finite
- the sampling design induces randomness, who is picked for the sample is random (not their response)

Design-based derivations

Use the design-based perspective to prove that \bar{y} is an unbiased estimator of the population mean \bar{y}_U when the design is a SRS.

$$\hat{\bar{y}}_U = \bar{y}$$

$$E(\bar{y}) = \bar{y}_U \text{ (pop. mean)}$$

↓
under SRS

Y_i 's fixed : Construct a random variable Z_i that indicates inclusion in the sample

$$Z_i = \begin{cases} 1, & \text{unit } i \text{ is included} \\ 0, & \text{otherwise} \end{cases}$$

$$\text{SRS: } Z_i \sim \text{Bern}(\pi_i = \frac{n}{N})$$

$$Z_i \sim \text{Bernoulli}(\pi_i)$$

For all units in pop.

$$Z_1, Z_2, \dots, Z_N$$

also:

$$\sum_{i=1}^N Z_i = n$$

Facts: $E(z_i) = 0(1 - \frac{n}{N}) + 1(\frac{n}{N}) = \frac{n}{N} = \pi_i$
 $Var(z_i) = (0 - \frac{n}{N})^2(1 - \frac{n}{N}) + (1 - \frac{n}{N})^2(\frac{n}{N}) = \frac{n}{N}(1 - \frac{n}{N}) = \pi_i(1 - \pi_i)$

Back to estimator

$$\underline{\bar{y}} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^N z_i y_i$$

$$E(\bar{y}) = E\left[\frac{1}{n} \sum_{i=1}^N z_i y_i\right] = \frac{1}{n} E\left[\sum_{i=1}^N z_i y_i\right] = \frac{1}{n} \sum_{i=1}^N \overbrace{E[z_i]}^{\frac{n}{N}} y_i$$

SRS $= \frac{1}{n} \sum_{i=1}^N \left(\frac{n}{N}\right) y_i = \frac{1}{N} \sum_{i=1}^N y_i = \bar{y}_u = \text{pop. mean}$

✓ $E(\bar{y}) = \bar{y}_u$ we know \bar{y} is unbiased!

Design-based derivations

Use the design-based perspective to prove that $\hat{t}_{HT} = \sum y_i / \pi_i$ is an unbiased estimator of the population total.

Prove: under any design with inclusion probs. π_i

$$E(\hat{t}_{HT}) = t = \sum_{i=1}^N y_i$$

Rewrite:

$$\hat{t}_{HT} = \sum_{i=1}^N \frac{y_i}{\pi_i} = \sum_{i=1}^N z_i \frac{y_i}{\pi_i}$$

$$\underline{E(\hat{t}_{HT})} = \sum_{i=1}^N \underbrace{E(z_i)}_{\pi_i} \frac{y_i}{\pi_i} = \sum_{i=1}^N \pi_i \cdot \frac{y_i}{\pi_i} = \sum_{i=1}^N y_i = t \quad \checkmark$$

so \hat{t}_{HT} is an unbiased estimator of t .

SRS variance

Prove that for a SRS:

$$\hat{y}_u = \bar{y} = \frac{1}{n} \sum_{i=1}^N z_i y_i$$

$$SE(\hat{y}_u) = \frac{SE(\hat{t})}{N} = \sqrt{\left(1 - \frac{n}{N}\right) \frac{s^2}{n}}$$

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^N z_i y_i\right) = (*)$$

If all z_i 's are independent, this is "easy" $\left(\frac{1}{n}\right)^2 \sum \text{Var}(z_i) y_i^2$

But when sampling without replacement from a finite population the z_i 's are dependent.

$$(*) = \left(\frac{1}{n}\right)^2 \sum_{i=1}^N \underbrace{\text{Var}(z_i)}_{\frac{n}{N} \left(1 - \frac{n}{N}\right)} y_i^2 + 2 \left(\frac{1}{n}\right)^2 \sum_{\substack{\text{all pairs} \\ (i,j)}} \text{Cov}(z_i, z_j) y_i y_j$$

What is Cov?

$$\text{Cov}(z_i, z_j) = E(z_i z_j) - E(z_i) E(z_j)$$

$\underbrace{\hspace{10em}}_{\text{only 1 when both = 1}} \quad \downarrow \frac{n}{N} \quad \downarrow \frac{n}{N}$

$$E(z_i z_j) = 1 \cdot \underbrace{P(z_i=1, z_j=1)}_{P(z_i=1)P(z_j=1|z_i=1)} + 0 \cdot \dots = \frac{n}{N} \times \frac{n-1}{N-1}$$

$$\begin{aligned} \text{Cov}(z_i, z_j) &= \frac{n}{N} \cdot \frac{n-1}{N-1} - \left(\frac{n}{N}\right)^2 = \frac{n}{N} \left(\frac{n-1}{N-1} - \frac{n}{N} \right) \\ &= \frac{n}{N} \left(\frac{N(n-1) - n(N-1)}{(N-1)N} \right) = \dots = \frac{n}{N} \left(1 - \frac{n}{N} \right) \left(\frac{-1}{N-1} \right) \end{aligned}$$

$$\begin{aligned} \text{SRS} \\ V(\bar{y}) &= \frac{1}{n^2} \sum y_i^2 \frac{n}{N} \left(1 - \frac{n}{N} \right) + 2 \frac{1}{n^2} \sum_{\text{all pairs}} y_i y_j \left(\frac{n}{N} \left(1 - \frac{n}{N} \right) \frac{-1}{N-1} \right) \text{negative!} \\ &= \text{lots of algebra} = \left(1 - \frac{n}{N} \right) \frac{s^2}{n} \end{aligned}$$