

# Sampling with unequal probabilities of selection (WOR)

Week 9 (6.4)

Stat 260, St. Clair

# Horvitz-Thompson Estimator

- You take a (one-stage) random sample where:
  - $t_i$  = "response" in PSU  $i$
  - $n$  = PSU sample size (or unique PSU sampled)
  - $\pi_i$  = sample inclusion prob for PSU  $i$
  - $w_i = 1/\pi_i$  = number of population units represented by PSU  $i$
- The Horvitz-Thompson estimator is

$$\hat{t}_{HT} = \sum_{i=1}^n w_i t_i$$

# Horvitz-Thompson Estimator

$$\hat{t}_{HT} = \sum_{i=1}^n w_i t_i$$

- For any design (with or without replacement), the H-T estimator is an unbiased estimator of population total  $t$ .

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

- Week 2 theory slides

# Horvitz-Thompson Estimator

- All total estimates so far, except for ratio estimates, have been HT estimators
  - SRS:  $w_i = N/n$
  - Stratified:  $w_{hj} = N_h/n_h$
  - One-stage cluster:  $w_{ij} = N/n$
- For these designs, the total estimator SE's can be derived from a general variance calculation
  - using the fact that **without replacement** designs leads to **dependence** among units being sampled

$$Var(\hat{t}_{HT}) = \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} t_i^2 + 2 \sum_{i=1}^N \sum_{\substack{k=1 \\ i < k}}^N \frac{\pi_{ik} - \pi_i \pi_k}{\pi_i \pi_k} t_i t_k$$

WOR  
 $\sum$  over all unordered pairs  $\{i, k\}$

# Horvitz-Thompson Estimator

$$Var(\hat{t}_{HT}) = \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} t_i^2 + 2 \sum_{\substack{i=1 \\ i < k}}^N \sum_{k=1}^N \frac{\pi_{ik} - \pi_i \pi_k}{\pi_i \pi_k} t_i t_k$$

- We need to compute the **joint** inclusion probability

$$\pi_{ik} = \pi_{ki} = P(\text{both } i, k \text{ included in the sample})$$

- What is  $\pi_{ik}$  for a SRS?

not order

$$\pi_{12} = P(\text{both 1 \& 2 in SRS})$$

sample that includes 1 & 2  $\{1, 2, \text{?}\}$   $n-2$  units to select

$$\pi_{12} = \frac{\text{\# of samples with units 1 \& 2}}{\binom{N}{n}} = \frac{\binom{N-2}{n-2}}{\binom{N}{n}} = \frac{n(n-1)}{N(N-1)}$$

# Horvitz-Thompson Estimator

- **SRS:** We measure  $t_i = y_i$  for each unit and  $\hat{t}_{HT} = N\bar{y}$ .

$$\pi_i = \frac{\binom{N-1}{n-1}}{\binom{N}{n}} = \frac{n}{N} \quad \pi_{ik} = \frac{\binom{N-2}{n-2}}{\binom{N}{n}} = \frac{n(n-1)}{N(N-1)}$$

- **SRS:** variance is then

$$\text{Var}(\hat{t}_{HT}) = \sum_{i=1}^N \frac{1 - \frac{n}{N}}{\frac{n}{N}} y_i^2 + 2 \sum_{i=1}^N \sum_{\substack{k=1 \\ i < k}}^N \frac{\frac{n(n-1)}{N(N-1)} - \frac{n}{N} \frac{n}{N}}{\frac{n}{N} \frac{n}{N}} y_i y_k$$

.... lots of algebra....

$$= N^2 \left(1 - \frac{n}{N}\right) \frac{S^2}{n} \Rightarrow \underline{\underline{\text{ch. 2}}}$$

theory

# Horvitz-Thompson Estimator

- All designs covered so far have used a SRS
  - Definition: each sample of size  $n$  is equally likely
  - Implication: each PSU is equally likely (equal  $\pi_i$  for all  $i$ )
- What if we don't use a SRS?
  - Take a random sample of PSU without replacement
  - Let inclusion probs  $\pi_i$  vary (unequal  $\pi_i$ )
- Our unbiased total estimate is still  $\hat{t}_{HT}$
- Our variance is computed using  $\pi_i$  and  $\pi_{ik}$ !

## Example: Unequal inclusion probabilities

- **Supermarkets** estimate total sales  $t$  for  $N = 4$  stores.

$t_i$  = total sales (thousands of dollars) at store  $i$

- **Design:** Random sample WOR with probability proportional to physical store size
- Here's our **population**:

Store	$t_i$	Size $m^2$
A	11	100
B	20	200
C	24	300
D	245	1000
total	$t = 300$	1600

↓  
what we want to estimate!



## Example: Unequal inclusion probabilities

**Design:** Random sample WOR with probability proportional to physical store size  $\rightarrow n = 2$  stores

$\psi_i$  = probability of selecting store  $i$  on your first draw

Store	$t_i$	Size $m^2$	$\psi_i$
A	11	100	$\psi_1 = \frac{100}{1600} = \frac{1}{16}$
B	20	200	$\psi_2 = 2/16$
C	24	300	$\psi_3 = 3/16$
D	245	1000	$\psi_4 = 10/16$
total	$t = 300$	1600	1

$$\psi_i = \frac{\text{size store } i}{\text{total of all store sizes}}$$

$\hookrightarrow 1600$

## Example: Unequal inclusion probabilities

**Design:** Random sample WOR with probability proportional to physical store size

$\psi_i$  = probability of selecting store  $i$  on your first draw

- **Catch:** WOR, probabilities for the second draw are not equal to  $\psi_i$ 
  - $\pi_i$  is the probability that store  $i$  is one of the two stores sampled, and not equal to  $\psi_i$

$\psi_1$  = prob. store 1 on first draw  
 $\neq \pi_1$  = prob. store 1 is picked in the sample of 2 stores.

## Example: Unequal inclusion probabilities

- Draw 1: we sample store B

$$\psi_{i|B} = P(\text{draw 2 is } i \mid \text{draw 1 is } B)$$

Store	Size $m^2$	$\psi_{i B}$
A	100	1/14
B	---	0
C	300	3/14
D	1000	10/14
total	1400	1

$= \psi_{A|B} = \frac{100}{1400} = \frac{1}{14}$   
 $\rightarrow \text{WOR}$   
 $= \psi_{C|B}$   
 $\psi_{D|B}$

## Example: Unequal inclusion probabilities

- Use the **individual PSU** selection probs (draw to draw) to compute the **joint inclusion** prob for each pair
- $\pi_{AB}$ : the probability that both A and B are included is

$$\begin{aligned} \downarrow \\ \pi_{AB} &= P(A \text{ then } B) + P(B \text{ then } A) \\ &= \psi_A \times \psi_{B|A} + \psi_B \psi_{A|B} \\ &\quad \downarrow \\ &\quad P(B \text{ on } 2^{\text{nd}} | A \text{ 1st}) \\ &= \frac{1}{16} \times \frac{200}{1500} \times \frac{2}{60} \times \frac{1}{14} \approx .0173 \end{aligned}$$

## Example: Unequal inclusion probabilities

- $\pi_A$ : the probability that store A is included is the **sum of all probs of samples that contain that PSU**

$$\begin{aligned} \pi_A &= \pi_{AB} + \pi_{AC} + \pi_{AD} \\ n=2 \quad &= .0173 + .0269 + .1458 \\ &= .19 \end{aligned}$$

## Example: Unequal inclusion probabilities

- The matrix below gives joint probs  $\pi_{ik}$  in the body and single PSU probs in the margins

	A	B	C	D	$\pi_i$
A	--	0.0173	0.0269	0.1458	0.1900 = $\pi_A$
B	0.0173	--	0.0556	0.2976	0.3705
C	0.0269	0.0556	--	0.4567	0.5393
D	0.1458	0.2976	0.4567	--	0.9002
$\pi_i$	0.1900	0.3705	0.5393	0.9002	$n = 2$

Handwritten notes:  $\pi_{AB}$  (above the cell for A, B),  $\pi_{BA}$  (to the left of the cell for B, A), and  $\pi_A$  (to the right of the cell for A,  $\pi_i$ ).

- Note that

theory

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

## Example: Unequal inclusion probabilities

- Suppose you sampled stores C and D

Store	Size $m^2$	$\pi_i$	$w_i$	$t_i$
C	300	0.5393	1.854	24
D	1000	0.9002	1.111	245

- Estimated total:

$$\hat{t}_{HT} = \sum_{\text{samped}} w_i t_i = (1.854)(24) + (1.111)(245) \\ = \$316.66$$

## Example: Unequal inclusion probabilities

- Variance of  $\hat{t}_{HT}$  is

$$\begin{aligned} \text{Var}(\hat{t}_{HT}) = & \left( \frac{1 - 0.1900}{0.1900} 11^2 + \dots + \frac{1 - 0.9002}{0.9002} 245^2 \right) \\ & + 2 \left( \frac{0.0173 - (0.1900)(0.3705)}{(0.1900)(0.3705)} \underset{A}{(11)} \underset{B}{(20)} + \right. \\ & \left. \dots + \frac{0.4567 - (0.5393)(0.9002)}{(0.5393)(0.9002)} \underset{C}{(24)} \underset{D}{(245)} \right) = 4383.6 \end{aligned}$$

- The estimated total sales is \$316.67 thousand with a SE of \$66.2 thousand.
- How does this compare to a SRS of  $n = 2$  stores?

$$\text{Var}(\hat{t}_{HT}) = \sum_{i=1}^4 \frac{1 - \pi_i}{\pi_i} t_i^2 + 2 \sum_{\substack{\text{all pairs} \\ \text{of store}}} \sum_k \frac{\pi_{ik} - \pi_i \pi_k}{\pi_i \pi_k} t_i t_k$$



## Example: Unequal inclusion probabilities

- Suppose stores C and D were selected from an SRS.
- **SRS Estimated total:** \$538 thousand

$$\hat{t}_{SRS} = N\bar{y} = 4 \frac{24 + 245}{2} = 538$$

- **SRS Variance** of  $\hat{t}_{SRS}$  is

$$Var(\hat{t}_{HT}) = 4^2 \left(1 - \frac{2}{4}\right) \frac{12874}{2} = 51496$$

*Handwritten notes: "ch. 2" with an arrow pointing to the variance term, and "S^2" with an arrow pointing to the 12874 term.*

$$\text{where } S^2 = \frac{1}{N-1} \sum_{i=1}^N (t_i - \bar{t}_U)^2 = 12874$$

```
> pop <- c(11,20,24,245)
> var(pop)
[1] 12874
```

## Example: Unequal inclusion probabilities

- **Probability proportional to size:** The estimated total sales is \$316.67 thousand with a SE of \$66.2 thousand.
- **SRS:** The estimated total sales is \$538 thousand with a SE of \$226.9 thousand.
- One important reason for selecting PSU with unequal probabilities:
  - can reduce SE (compared to SRS) when selection probability  $\pi_i$  is positively associated with the response  $t_i$
  - called probability proportional to size (pps) sampling  $\hat{\pi}_i \propto t_i$
  - most samples will contain large  $t_i$  making variation in  $\hat{t}_{pps}$  less than when a small  $t_i$  is just as likely as a large

Store: probs. prop. to store size.  
→ size &  $t_i$  (sales) are positively correlated

## Example: Unequal inclusion probabilities

- One important reason for **not** selecting PSU with unequal probabilities:
  - if some PSU have very small  $\pi_i$ , then they have very high weight  $w_i$
  - can cause imprecise estimates of  $Var(\hat{t}_{HT})$  because of these high weights

$$\frac{1}{\pi_i} \approx w_i$$

very small  $\swarrow$   $\searrow$  very big

# Estimating HT variance

$$Var(\hat{t}_{HT}) = \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} t_i^2 + 2 \sum_{i=1}^N \sum_{\substack{k=1 \\ i < k}}^N \frac{\pi_{ik} - \pi_i \pi_k}{\pi_i \pi_k} t_i t_k$$

- To estimate variance, treat the summations as population totals
  - estimate the total with a HT-estimator!
  - weight sampled values by  $w_i$ 's

theory

# Estimating HT variance

- There are three commonly used estimates of  $Var(\hat{t}_{HT})$
- **Horvitz-Thompson (HT)**: unbiased and often the default software version (as in survey), but can be negative for samples with small inclusion probs

$$\hat{V}_{HT}(\hat{t}_{HT}) = \sum_{i=1}^n \frac{1 - \pi_i}{\pi_i^2} t_i^2 + 2 \sum_i \sum_{\substack{k \\ i < k}} \frac{\pi_{ik} - \pi_i \pi_k}{\pi_{ik}} \frac{t_i}{\pi_i} \frac{t_k}{\pi_k}$$

sample unit

pairs

- **Sen-Yates-Grundy (SYG)**: unbiased and more stable than HT version

$$\hat{V}_{SYG}(\hat{t}_{HT}) = \sum_i \sum_{\substack{k \\ i < k}} \frac{\pi_i \pi_k - \pi_{ik}}{\pi_{ik}} \left( \frac{t_i}{\pi_i} - \frac{t_k}{\pi_k} \right)^2$$

sample units

# Estimating HT variance

- **With Replacement:** is a biased estimate that overestimates the variance, but it doesn't require joint inclusion probs!

$$\hat{V}_{WR}(\hat{t}_{HT}) = \frac{n}{n-1} \sum_{i=1}^n \left( \frac{t_i}{\pi_i} - \frac{\hat{t}_{HT}}{n} \right)^2$$

---

$\hat{V}_{HT}(\hat{t}_{HT}) \rightarrow$  stores C&D sample  $\pi_c = .5393$   $\pi_D = .9002$   
 $t_c = 24$   $t_D = 245$   $\pi_{CD} = .4567$

$$\begin{aligned} \hat{V}_{HT}(\hat{t}_{HT}) &= \frac{1 - .5393}{.5393^2} (24)^2 + \frac{1 - .9002}{.9002^2} (245)^2 \\ &+ 2 \times \left( \frac{.4567 - (.5393)(.9002)}{.4567} \right) \frac{24}{.5393} \times \frac{245}{.9002} \\ &\approx 6778 \end{aligned}$$

## Example: Estimating HT variance

→  $\pi_i K$

Sample $\mathcal{S}$	$P(\mathcal{S})$	$\hat{t}_{HT}$	$\hat{V}_{HT}(\hat{t}_{HT})$	$\hat{V}_{SYG}(\hat{t}_{HT})$
A,B	0.01726	111.87	-14,691.5	47.1
A,C	0.02692	102.39	-10,832.1	502.8
A,D	0.14583	330.06	4,659.3	7,939.8
B,C	0.05563	98.48	-9,705.1	232.7
B,D	0.29762	326.15	5,682.8	5,744.1
C,D	0.45673	316.67	<u>6,782.8</u>	3,259.8

- If we happen to sample two small stores, our HT estimate of variance is negative!
- But both are unbiased estimators of the true variance.

$$E[\hat{V}_{HT}] = (-14691.5)(.01726) + \dots = V_{HT} = 4383$$

↪ all samples

$$E[\hat{V}_{SYG}] = 47.1(.01726) + \dots = V_{HT} = 4383$$

# What about estimating population mean?

$$\sum_{\text{all elements}} w_i$$

→ observation unit

- Summing sampling weights over all **elements** sampled will give
  - actual population size (of elements) when weights are equal for all elements and number of elements per PSU is constant
  - an unbiased estimated population size (of elements)
- The Horvitz-Thompson estimate of population mean (per element) is

$$\hat{y}_{HT} = \frac{\hat{t}_{HT}}{\sum_{\text{all elements}} w_i} = \frac{\sum w_i t_i}{\sum w_i}$$

- The survey package uses this when you run `svymean`
  - gives  $\hat{t}_{HT}$  when you run `svytotal`

