

Sampling with unequal probabilities of selection (WOR) - theory

Week 9 (6.4)

Stat 260, St. Clair

Horvitz-Thompson Estimator

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

Derive the variance of \hat{t}_{HT} :

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

theory
ch. 2

$$Z_i = \begin{cases} 1, & \text{if } i \text{ in sample} \\ 0, & \text{not in sample} \end{cases} \quad Z_i \sim \text{Bern}(\pi_i)$$

$$P(Z_i = 1) = P(\text{unit } i \text{ in sample}) = \pi_i$$

$$E(Z_i) = \pi_i \quad Var(Z_i) = \pi_i (1 - \pi_i)$$

$$E[Z_i Z_k] = 1 \times P(\text{units } i + k \text{ in sample}) + 0 \cdot \text{---} = \pi_{ik}$$

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

ch. 2 thg $\rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$

$$\text{Var}(\hat{t}_{HT}) = \text{Var}\left(\sum_{i=1}^N z_i \frac{t_i}{\pi_i}\right)$$

$$= \underbrace{\sum_{i=1}^N \text{Var}\left(z_i \frac{t_i}{\pi_i}\right)}_{\text{all pairs } i, k} + 2 \underbrace{\sum_{i=1}^N \sum_{\substack{k=1 \\ i < k}}^N \text{Cov}\left(z_i \frac{t_i}{\pi_i}, z_k \frac{t_k}{\pi_k}\right)}_{\text{all pairs } i, k}$$

$$= \sum_{i=1}^N \left(\frac{t_i}{\pi_i}\right)^2 \text{Var}(z_i) + 2 \sum_{\text{all pairs}} \sum \frac{t_i}{\pi_i} \frac{t_k}{\pi_k} \underbrace{\text{Cov}(z_i, z_k)}_{\uparrow}$$

$$\text{Cov}(z_i, z_k) = E[z_i z_k] - E[z_i] E[z_k] = \pi_{ik} - \pi_i \pi_k$$

$$\text{Var}(\hat{t}_{HT}) = \sum_{i=1}^N \left(\frac{t_i}{\pi_i} \right)^2 \pi_i (1 - \pi_i) + 2 \sum_i \sum_k \frac{t_i t_k}{\pi_i \pi_k} (\pi_{ik} - \pi_i \pi_k)$$

\downarrow
 Var.

$$= \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} t_i^2 + 2 \sum_i \sum_k \frac{\pi_{ik} - \pi_i \pi_k}{\pi_i \pi_k} t_i t_k$$

pairs(i, k)



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

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

.....lots of algebra.....

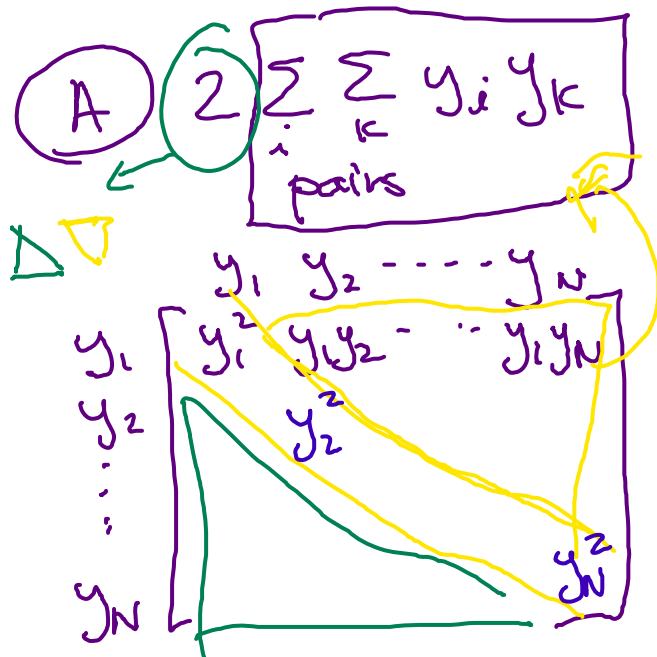
$$= N^2 \left(1 - \frac{n}{N}\right) \frac{S^2}{n}$$

Fill in the missing work...

→ ch. 2

$$\text{Var}(\hat{t}_{HT}) = \frac{N}{n} \left(1 - \frac{n}{N}\right) \underbrace{\sum_{i=1}^N y_i^2}_B + \underbrace{\frac{N}{n} \left[\frac{n-1}{N-1} - \frac{n}{N} \right] 2 \sum_{\substack{\hat{i}=\hat{k} \\ \text{pairs}}} y_i y_k}_A$$

$$\begin{aligned} \textcircled{B} \quad \frac{N}{n} \left[\frac{n-1}{N-1} \cdot \frac{N}{N} - \frac{n}{N} \cdot \frac{N-1}{N-1} \right] &= \frac{N}{n} \left[\frac{n(n-1) - n(N-1)}{N(N-1)} \right] \\ &= \frac{N}{n} \left[\frac{Nn - N - nN + n}{N(N-1)} \right] = \frac{N}{n} \left[\frac{n - N}{N(N-1)} \right] = \frac{N}{n} \cdot \frac{-1}{N-1} \left(1 - \frac{n}{N}\right) \end{aligned}$$



$$\begin{aligned} &= \frac{N^2}{N^2} \sum_{i=1}^N \sum_{k=1}^N y_i y_k - \sum_{i=1}^N y_i^2 = N^2 \sum_{i=1}^N \frac{y_i}{N} \sum_{k=1}^N \frac{y_k}{N} - \sum_{i=1}^N y_i^2 \\ &= N^2 \bar{y}^2 - \sum_{i=1}^N y_i^2 \end{aligned}$$

SRS
↓

(A)
↓

$$\text{Var}(\hat{\tau}_{HT})$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \left[\sum_{i=1}^N y_i^2 \times \frac{N-1}{N-1} - \frac{1}{N-1} \left(N^2 \bar{y}_u^2 - \sum y_i^2 \right) \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[N \sum_{i=1}^N y_i^2 - \sum_{i=1}^N y_i^2 - N^2 \bar{y}_u^2 + \sum_{i=1}^N y_i^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \left(\frac{N}{N-1} \right) \left[\sum y_i^2 - N \bar{y}_u^2 \right]$$

ch. 2 thm = $\sum_{i=1}^N (y_i - \bar{y}_u)^2$

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y}_u)^2$$

$$\begin{aligned}\text{Var}(\hat{t}_{HT}) &= \text{Var}(N \bar{y}) \\ &= N^2 \left(1 - \frac{n}{N}\right) \frac{S^2}{n}\end{aligned} \quad \checkmark$$

Estimating HT variance

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

Derive the HT-estimator of $Var(\hat{t}_{HT})$

proposed way

$$\hat{var}_{HT}(\hat{t}) = \sum_{i=1}^N w_i z_i \textcircled{A} + 2 \sum \sum_{\text{pairs}} w_{ik} z_i z_k \textcircled{B}$$

$\frac{1}{\pi_{ik}} = E(z_i z_k)$

$$= \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i \cdot \pi_i} t_i^2 + 2 \sum_{\substack{i \\ \text{obs.}}} \sum_k \frac{\pi_{ik} - \pi_i \pi_k}{\pi_{ik}} \cdot \frac{t_i}{\pi_i} \cdot \frac{t_k}{\pi_k}$$

HT - variance est.

Horvitz-Thompson Estimator

Prove that

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

→ use fact $E[z_i] = \pi_i$

$$\underline{\underline{\sum_{i=1}^N \pi_i}} = \sum_{i=1}^N E[z_i] = E\left[\sum_{i=1}^N z_i\right] = E[n] = \underline{\underline{n}}$$

↓
must equal n!

Horvitz-Thompson Estimator

$$\frac{\sum w_i t_i}{\sum w_i}$$

Prove that if there are M observation units in the population,

$$E\left(\sum_{\text{all sampled obs. units}} w_i\right) = M$$

$$E\left[\sum_{\substack{\text{all} \\ \text{sampled} \\ \text{obs. units}}} w_i\right] = E\left[\sum_{\substack{\text{all pop.} \\ \text{obs.} \\ \text{unit}}} z_i w_i\right] = \sum_{\substack{\text{all pop.} \\ \text{obs. unit}}} w_i E[z_i]$$

$\nearrow \frac{1}{\pi_i}$

$$= \sum_{\substack{\text{all pop.} \\ \text{obs. unit}}} \left(\frac{1}{\pi_i}\right) \pi_i = \sum_{\substack{\text{all pop.} \\ \text{obs. unit}}} 1 = \underbrace{(1 + 1 + \dots + 1)}_{M \text{ obs. units}} = \underline{\underline{M}}$$