

Stat 576  
5-13-25

$$\text{Also need } E\left[\frac{N}{n} \sum_{i=1}^n M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right)\right] \quad (1)$$

$$= E_1 E_2 L \cdot 1$$

$$= N E_1 E_2 \left[ \frac{1}{n} \sum_{i=1}^n M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right) \right]$$

$$= N E_1 \left[ \frac{1}{n} \sum_{i=1}^n M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right) \right]$$

$$= N \frac{1}{N} \sum_{i=1}^N M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right)$$

$$E(\hat{V}[\hat{\epsilon}]) =$$

$$\frac{N^2}{n} \left(1 - \frac{n}{N}\right) \left[ \frac{1}{N} \sum_{i=1}^N M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right) + s_e^2 \right]$$

$$+ \sum_{i=1}^N M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right)$$

$$= \frac{N^2}{n} s_e^2 \left(1 - \frac{n}{N}\right) + \frac{N}{n} \sum_{i=1}^N M_i^2 \frac{s_i^2}{m_i} \left(1 - \frac{m_i}{M_i}\right)$$

$\therefore \hat{V}[\hat{\epsilon}]$  is actually  
an unbiased estimator of  $V[\hat{\epsilon}]$

$$\uparrow \frac{N}{n} \left(1 - \frac{n}{N}\right) + 1$$

(3)

Single-stage cluster sampling  
with unequal probabilities

Assume we sample  $n$  clusters out of  $N$ , WR.

Let  $z_i$  be the probability that cluster  $i$   
is selected on a particular draw.

$$\text{Let } \hat{t} = \frac{1}{n} \sum_{i=1}^n \frac{t_i}{z_i}$$

$$\text{Then } E[\hat{t}] = \frac{1}{n} \sum_{i=1}^n E\left[\frac{t_i}{z_i}\right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^N \frac{t_j}{z_j} \cdot z_i \right)$$

$$= \frac{1}{n} \sum_{i=1}^n t = t$$

$$\text{Also, } V[\hat{t}] = V\left[ \frac{1}{n} \sum_{i=1}^n \frac{t_i}{z_i} \right]$$

$$= \frac{1}{n^2} \sum_{i=1}^n V\left[\frac{t_i}{z_i}\right]$$

$$= \frac{1}{n^2} \sum_{i=1}^n \underbrace{\left( \sum_{j=1}^N \left( \frac{t_j}{z_j} - t \right)^2 \cdot z_i \right)}_{\sigma_z^2} = \frac{\sigma_z^2}{n}$$

(4)

(5)

Note: if you could set  $\bar{z}_j = \frac{t_j}{t} \bar{y}_j$ ,

then  $\sigma_z^2 = 0$  and your estimator would be perfect.

Since the  $t_j$ 's are unknown (as is  $t$ ).

One possibility is to use  $\hat{\psi}_j = \frac{M_j}{K}$

This is called PPS sampling

↑ probabilities proportional  
to size

(6)

$$\hat{t}_{PPS} = \frac{1}{n} \sum_{i=1}^n \frac{t_i}{\psi_i} = \frac{1}{n} \sum_{i=1}^n \frac{t_i}{M_i/K}$$

$$= \frac{K}{n} \sum_{i=1}^n \bar{y}_i = K \bar{y}$$

$$V[\hat{t}_{PPS}] = \frac{1}{n} \sigma_z^2 = \frac{1}{n} \underbrace{\sum_{j=1}^n \left( \frac{t_j}{\psi_j} - \bar{t} \right)^2 \psi_j}_{\sigma_{PPS}^2}$$

Estimate  $\sigma_{PPS}^2$  with  $s_{PPS}^2$ .

$s_{PPS}^2$  is the sample variance of the  $\frac{t_j}{\psi_j}$  terms

(7)

$$\text{But } \frac{\bar{t}_i}{\psi_i} = K \bar{y}_i$$

So  $s_{\text{pps}}^2$  is the sample variance of the  $K \bar{y}_i$  terms,  
which is  $K^2 s_{\bar{y}}^2$

$\uparrow$  Sample variance of  
the cluster means

Summary:

$$\hat{t}_{\text{pps}} = K \bar{y}$$

$$V[\hat{t}_{\text{pps}}] = \frac{\sigma_{\text{pps}}^2}{n}$$

$$\hat{V}[\hat{t}_{\text{pps}}] = \frac{K^2 s_{\bar{y}}^2}{n}$$

$$\bar{y}_{\text{pps}} = \bar{y}$$

$$V[\bar{y}_{\text{pps}}] = \frac{\sigma_{\text{pps}}^2}{n K^2}$$

$$\hat{V}[\bar{y}_{\text{pps}}] = \frac{s_{\bar{y}}^2}{n}$$

(8)

A generalization of the idea behind  
PPS sampling

let  $\pi_j$  be the probability that  $y_j$  is in  
the sample.

let  $\pi_{ij}$  be the joint probability that  $y_i$  and  $y_j$   
are both in the sample

Let  $\hat{t}_{HT} = \sum_{i=1}^n \frac{y_i}{\pi_i}$  Horvitz-Thompson  
estimator

(9)

Hörvitz-Thompson theorem:

$$E(\hat{t}_{HT}) = t$$

$$\text{Proof: } E[\hat{t}_{HT}] = E\left[\sum_{i=1}^N \frac{Y_i}{\pi_i}\right]$$

$$= E\left[\sum_{j=1}^N \frac{Y_j}{\pi_j} U_j\right] \quad \text{let } U_j = \begin{cases} 1 & \text{if } Y_j \text{ is in sample} \\ 0 & \text{otherwise} \end{cases}$$

$$= \sum_{j=1}^N \frac{Y_j}{\pi_j} E(U_j) \quad U_j \sim \text{Bernoulli}(\pi_j)$$

$$= \sum_{j=1}^N \frac{Y_j}{\pi_j} \pi_j = t$$