

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

$$= E_1 E_2 [\cdot]$$

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

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

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

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

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

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

So $\hat{V}[\hat{\tau}]$ is actually $\uparrow \frac{N}{n} \left(1 - \frac{n}{N}\right) + 1$
 an unbiased estimator of $V(\hat{\tau})$

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.

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

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

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

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

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

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

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

Note: if you could set $z_j = \frac{t_j}{t} \psi_j$,

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 $\psi_j = \frac{M_j}{K}$

This is called PPS sampling

↑ probabilities proportional to size

$$\begin{aligned}\hat{t}_{pps} &= \frac{1}{n} \sum_{i=1}^n \frac{t_i}{\psi_i} = \frac{1}{n} \sum_{i=1}^n \frac{t_i}{M_i/K} \\ &= \frac{K}{n} \sum_{i=1}^n y_i = K \bar{y}\end{aligned}$$

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

Estimate σ_{pps}^2 with s_{pps}^2 .

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

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

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

↑ sample variance of
the cluster means

Summary:

$$\begin{aligned}\hat{t}_{pps} &= K \bar{y} \\ V[\hat{t}_{pps}] &= \frac{\sigma_{pps}^2}{n} \\ \hat{V}[\hat{t}_{pps}] &= \frac{K^2 s_{\bar{y}}^2}{n}\end{aligned}$$

$$\begin{aligned}\bar{y}_{pps} &= \bar{y} \\ V[\bar{y}_{pps}] &= \frac{\sigma_{pps}^2}{n K^2} \\ \hat{V}[\bar{y}_{pps}] &= \frac{s_{\bar{y}}^2}{n}\end{aligned}$$

Stat 576 HW#6

- 9** The file `statepps.dat` lists the number of counties, land area, and 1992 population for the 50 states plus the District of Columbia.
- a** Use the cumulative-size method to draw a sample of size 10 with replacement, with probabilities proportional to land area. What is ψ_i for each state in your sample?
 - b** Use the cumulative-size method to draw a sample of size 10 with replacement, with probabilities proportional to population. What is ψ_i for each state in your sample?
 - c** How do the two samples differ? Which states tend to be in each sample?