

### References

- 1. Cochran, W.G. (1963) Sampling Techniques, Wiley, New York.
- 2. Kish, L. (1995) Survey Sampling, Wiley Inter. Science.
- 3. Lohr, S.L. (1999) Sampling: Design and Analysis, Duxbury Press.
- 4. McLennan, W. (1999) An Introduction to Sample Surveys, A.B.S. Publications, Canberra.

# Section outline

- Simple random samples and stratification.
   Finite population correction factor. Sample size determination. Inference over subpopulations.
- 2. Stratified sampling. Optimal allocation.
- Ratio and regression estimators. Ratio estimators. Hartley-Ross estimator. Ratio estimator for stratified samples. Regression estimator.
- 4. Systematic sampling and cluster sampling.
- Sampling with unequal probabilities. Probability proportional to size(PPS) sampling. The Horvitz-Thompson estimator.



# 1 Simple Random Samples (SRS)

## 1.1 The Population

We have a finite number of elements, N where N is assumed known. The population is  $Y_1 \ldots Y_N$ , where  $Y_i$  is a numerical value associated with *i*-th element. We adopt the notation where capital letters refer to characteristics of the population; small letters are used for the corresponding characteristics of a sample.

Population Total:  $Y = \sum_{i=1}^{N} Y_i$ , Population Mean:  $\mu = \bar{Y} = \frac{Y}{N} = \frac{1}{N} \sum_{i=1}^{N} Y_i$ , Population Variance :  $\sigma^2 = \frac{1}{N} \sum_{i=1}^{N} (Y_i - \bar{Y})^2$  and  $S^2 = \frac{N}{N-1} \sigma^2 = \frac{1}{N-1} \sum_{i=1}^{N} (Y_i - \bar{Y})^2$ .

These are fixed (population) quantities, to be estimated.

If we have to consider two numerical values:  $(Y_i, X_i)$ ,  $i = 1, \dots, N$ an additional population quantity of interest is

$$R = \frac{\sum\limits_{i=1}^{N} Y_i}{\sum\limits_{i=1}^{N} X_i} = \frac{Y}{X} = \frac{\bar{Y}}{\bar{X}}$$

the ratio of totals.

THE UNIVERSITY OF SYDNEY

### 1.2 Simple Random Sampling

Focus on the numerical values  $Y_i$ ,  $i = 1, \dots, N$ . A random sample of size n is taken without replacement: the observed values  $y_1, \dots, y_n$ are random variables and are stochastically dependent. The sampling frame is a list of the values  $Y_i$ ,  $i = 1, \dots, N$ .

The natural estimator for 
$$\overline{Y}$$
 is  $\widehat{\overline{Y}} = \frac{1}{n} \sum_{i=1}^{n} y_i = \overline{y}$   
and hence for  $Y = N\mu$  is  $\widehat{Y} = N\overline{y}$ .

Distributional properties of  $\bar{y}$  are complicated by the dependence of the  $y_i$ 's.

Sample variance: 
$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (y_{i} - \bar{y})^{2}$$

### **Fundamental Results**

$$E(\bar{y}) = \mu.$$
  

$$\operatorname{Var}(\bar{y}) = \left(1 - \frac{n}{N}\right)\frac{S^2}{n} = (1 - f)\frac{S^2}{n} = \left(\frac{N - n}{N - 1}\right)\frac{\sigma^2}{n}$$
  

$$\operatorname{var}(\bar{y}) = \left(1 - \frac{n}{N}\right)\frac{s^2}{n}$$
  

$$E(s^2) = S^2$$

where f is the sampling fraction and the finite population correction (f.p.c.) is 1 - f.



**Proof:** Let 
$$\overline{y} = \frac{1}{n} \sum_{i \in S} y_i = \frac{1}{n} \sum_{i=1}^N y_i I_i$$
 where the sample membership

indicator

$$I_i = \begin{cases} 1 & \text{if element } i \text{ is in the sample,} \\ 0 & \text{if otherwise.} \end{cases}$$

First, we have

$$E(I_i) = 0 \times \Pr(I_i = 0) + 1 \times \Pr(I_i = 1) = \pi_i = \frac{n}{N},$$

$$E(I_i^2) = 0^2 \times \Pr(I_i = 0) + 1^2 \times \Pr(I_i = 1) = \pi_i = \frac{n}{N},$$

$$E(I_iI_j) = 0 \times 0 \Pr(I_i = 0 \& I_j = 0) + 0 \times 1 \Pr(I_i = 0 \& I_j = 1) + 1 \times 0 \Pr(I_i = 1 \& I_j = 0) + 1 \times 1 \Pr(I_i = 1 \& I_j = 1)$$

$$= \pi_{ij} = \frac{n(n-1)}{N(N-1)}$$

$$\operatorname{Var}(I_i) = E(I_i^2) - E^2(I_i) = \pi_i(1 - \pi_i) = \frac{n}{N}(1 - \frac{n}{N}),$$

$$\operatorname{Cov}(I_i, I_j) = E(I_iI_j) - E(I_i)E(I_j) = \pi_{ij} - \pi_i\pi_j = \frac{n(n-1)}{N(N-1)} - \left(\frac{n}{N}\right)^2.$$

Then

$$E(\overline{y}) = \frac{1}{n} \sum_{i=1}^{N} y_i E(I_i) = \frac{1}{n} \sum_{i=1}^{N} y_i \cdot \frac{n}{N} = \frac{1}{N} \sum_{i=1}^{N} y_i = \overline{Y}. \quad \text{Unbiased}$$
$$E[\operatorname{var}(\overline{y})] \stackrel{def.}{=} E\left[\left(\frac{N-n}{N}\right) \frac{s_y^2}{n}\right] \stackrel{Pf.2 \to}{=} \left(\frac{N-n}{N}\right) \frac{S_y^2}{n} \stackrel{\leftarrow Pf.1}{=} \operatorname{Var}(\overline{y}) \text{ Unbiased}$$

**Proof 1:** show that  $\operatorname{Var}(\overline{y}) = \left(\frac{N-n}{N}\right) \frac{S_y^2}{n}$ .



$$\begin{aligned} \operatorname{Var}(\overline{y}) &= \operatorname{Var}\left(\frac{1}{n}\sum_{i=1}^{N} y_{i}I_{i}\right) \\ &= \frac{1}{n^{2}} \left[\sum_{i=1}^{N} y_{i}^{2}\operatorname{Var}(I_{i}) + 2\sum_{i}\sum_{j,i < j} y_{i}y_{j}\operatorname{Cov}(I_{i}, I_{j})\right] \\ &= \frac{1}{n^{2}} \left\{\sum_{i=1}^{N} y_{i}^{2} \left[\frac{n}{N}\left(1 - \frac{n}{N}\right)\right] + 2\sum_{i}\sum_{j,i < j} y_{i}y_{j} \left[\frac{n(n-1)}{N(N-1)} - \left(\frac{n}{N}\right)^{2}\right]\right\} \\ &= \frac{n}{n^{2}} \left\{\frac{1}{N}\left(1 - \frac{n}{N}\right)\sum_{i=1}^{N} y_{i}^{2} + 2\frac{1}{N}\left(\frac{n-1}{N-1} - \frac{n}{N}\right)\sum_{i}\sum_{j,i < j} y_{i}y_{j}\right\} \\ &= \frac{1}{n}\left(1 - \frac{n}{N}\right) \left\{\frac{1}{N}\sum_{i=1}^{N} y_{i}^{2} + 2\frac{1}{N}\left(1 - \frac{n}{N}\right)^{-1}\frac{N(n-1) - n(N-1)}{N(N-1)}\sum_{i}\sum_{j,i < j} y_{i}y_{j}\right\} \\ &= \frac{1}{n}\left(1 - \frac{n}{N}\right) \left\{\frac{1}{N}\sum_{i=1}^{N} y_{i}^{2} + 2\frac{1}{N}\frac{N}{N-n}\frac{n-N}{N(N-1)}\sum_{i}\sum_{j,i < j} y_{i}y_{j}\right\} \\ &= \frac{1}{n}\left(1 - \frac{n}{N}\right)\frac{1}{N-1} \left\{\sum_{i=1}^{N} y_{i}^{2} - \frac{1}{N}\left(\sum_{i=1}^{N} y_{i}^{2} + 2\sum_{i}\sum_{j,i < j} y_{i}y_{j}\right)\right\} \\ &= \frac{1}{n}\left(1 - \frac{n}{N}\right)\frac{1}{N-1} \left\{\sum_{i=1}^{N} y_{i}^{2} - \frac{1}{N}\left(\sum_{i=1}^{N} y_{i}^{2} + 2\sum_{i}\sum_{j,i < j} y_{i}y_{j}\right)\right\} \\ &= \frac{1}{n}\left(1 - \frac{n}{N}\right)\frac{1}{N-1} \left\{\sum_{i=1}^{N} y_{i}^{2} - \frac{1}{N}\left(\sum_{i=1}^{N} y_{i}\right)\left(\sum_{i=1}^{N} y_{i}\right)\right\} \\ &= \frac{1}{n}\left(1 - \frac{n}{N}\right)\frac{1}{N-1} \left\{\sum_{i=1}^{N} y_{i}^{2} - \frac{1}{N}\left(\sum_{i=1}^{N} y_{i}\right)\left(\sum_{i=1}^{N} y_{i}\right)\right\} \\ &= \left(1 - \frac{n}{N}\right)\frac{S_{y}^{2}}{n} = \frac{N-n}{N}\frac{N}{N-1}\frac{\sigma_{y}^{2}}{n} = \left(\frac{N-n}{N-1}\right)\frac{\sigma_{y}^{2}}{n} \end{aligned}$$

**Proof 2:** show that  $E(s_y^2) = S_y^2$  where  $S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \overline{Y})^2 = \frac{N}{N-1} \sigma_y^2$ .

$$E(s_y^2) = E\left[\frac{1}{n-1}\sum_{i=1}^n (y_i - \overline{y})^2\right] = \frac{1}{n-1}E\left\{\sum_{i=1}^n [(y_i - \overline{Y}) - (\overline{y} - \overline{Y})]^2\right\}$$



$$= \frac{1}{n-1} E\left[\sum_{i=1}^{n} (y_i - \overline{Y})^2 - n(\overline{y} - \overline{Y})^2\right]$$

$$= \frac{1}{n-1} \left[\sum_{i=1}^{n} E(y_i - \overline{Y})^2 - nE(\overline{y} - \overline{Y})^2\right]$$

$$= \frac{1}{n-1} \left[\sum_{i=1}^{n} Var(y_i) - nVar(\overline{y})\right]$$

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

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

$$= \frac{N}{N-1}\sigma_y^2 = S_y^2$$

### **Central Limit Property**

For large sample size n (n > 30 say), and small to moderate f, we have the approximation

$$(\bar{y} - \mu) / \sqrt{\operatorname{Var}(\bar{y})} \sim \mathcal{N}(0, 1).$$

# Confidence Interval for $\mu = \bar{Y}$

Replacing  $S^2$  by  $s^2$ , an approximate 95% C.I. for  $\mu$  and Y are respectively

$$\bar{y} \pm 1.96 \frac{s}{\sqrt{n}} \sqrt{1-f},$$
  
 $N(\bar{y} \pm 1.96 \frac{s}{\sqrt{n}} \sqrt{1-f})$ 

Read Tutorial 10 Q1.



**Example:** (Industrial firm) An industrial firm is concerned about the time spent each week by staff on certain tasks. The time-log sheets of a SRS of n = 50 employees show the average amount of time spent on these tasks is 10.31 hours, with a sample variance  $s^2 = 2.25$ . The company employs N = 750 staff. Estimate the total number of man-hours used each week on the tasks and construct a 95% CI for the estimate.

**Solution:** From N = 750 time-log sheets, a SRS of n = 50 sheets was obtained. The average amount of time used in the sample is  $\overline{y} = 10.31$  hours/week. Since n = 50 is large, we take  $z_{0.025} = 1.96$ . Hence

### 1.3 Simple Random Sampling for Attributes.

The method for SRS can be applied to estimate the *total number*, or *proportion* (or %) of units which possess some qualitative attribute. Let this subset of the population be C.

Let

$$Y_i = 1 \quad \text{if} \quad i \in C$$
$$= 0 \quad \text{if} \quad i \notin C$$

and similarly for  $y_i$ 's .

Customary notation:

 $\mu = \bar{Y} = P \text{ is the population proportion,}$   $Y = \sum_{i=1}^{N} Y_i = NP \text{ is the population total count,}$  $\bar{y} = p \text{ is the sample proportion.}$ 

Let 
$$y = \sum_{i=1}^{n} y_i = np$$
,  $Q = 1 - P$  and  $q = 1 - p$ .  
Thus  $\widehat{Y} = Np = Ny/n$ .

Since

$$\sum_{i=1}^{N} (Y_i - \bar{Y})^2 = \sum_{i=1}^{N} Y_i^2 - N\bar{Y}^2 = N\bar{Y} - N\bar{Y}^2 = NP(1-P),$$

we have

$$S^{2} = NP(1-P)/(N-1) \approx P(1-P)$$

and similarly,

$$s^2 = np(1-p)/(n-1) \approx p(1-p)$$
 and  
 $\frac{s^2}{n} = \frac{np(1-p)}{n(n-1)} = \frac{p(1-p)}{n-1} \approx \frac{p(1-p)}{n}$ 

#### Sample Size Calculations 1.4

To calculate the sample size needed for sampling yet to be carried out, we want to be at least  $100(1-\alpha)\%$  sure the estimate  $\bar{y}$  of  $\mu$  is within 1006% of the actual value of  $\mu$  (e.g.  $1 - \alpha = 0.95$ ,  $z_{\alpha/2} = 1.96$ ). That is

$$\Pr\{|\bar{y} - \mu| \leq \delta_{\mu}\} \geq 1 - \alpha$$

$$\Leftrightarrow \Pr\left(\frac{|\bar{y} - \mu|}{\sqrt{\operatorname{Var}(\bar{y})}} \leq \frac{\delta_{\mu}}{\sqrt{\operatorname{Var}(\bar{y})}}\right) \geq 1 - \alpha$$

$$\Leftrightarrow \delta_{\mu} / \sqrt{\operatorname{Var}(\bar{y})} \geq z_{\alpha/2}$$

$$\Leftrightarrow (1 - \frac{n}{N}) \frac{S^{2}}{n} \leq \frac{(\delta_{\mu})^{2}}{z_{\alpha/2}^{2}}$$

$$\Leftrightarrow \frac{S^{2}}{n} - \frac{S^{2}}{N} \leq \frac{(\delta_{\mu})^{2}}{z_{\alpha/2}^{2}}$$

$$\Leftrightarrow \frac{S^{2}}{n} \leq \frac{(\delta_{\mu})^{2}}{z_{\alpha/2}^{2}} + \frac{S^{2}}{N}$$

$$\Leftrightarrow n \geq \frac{S^{2}}{\frac{(\delta_{\mu})^{2}}{z_{\alpha/2}^{2}} + \frac{S^{2}}{N}}$$
Normal 
$$\lim_{\substack{\lambda = 1 - \alpha \\ \frac{-\delta_{\mu}}{SE(\bar{y})} = 0}} \frac{1 - \alpha}{z_{\alpha/2}} = 1.96$$

$$\Leftrightarrow n \geq \frac{NS^{2}}{N(\delta_{\mu})^{2}/z_{\alpha/2}^{2} + S^{2}}$$

Ignoring f.p.c. (i.e. taking f = 0 or  $\operatorname{Var}(\bar{y}) = S^2/n$  when N is unknown)  $\boxed{n \ge \frac{z_{\alpha/2}^2 S^2}{(S_{\alpha})^2} \approx \frac{z_{\alpha/2}^2 s^2}{(S_{\alpha})^2}}$ 

$$n \ge rac{z_{lpha/2}^2 S^2}{(\delta_\mu)^2} \, pprox \, rac{z_{lpha/2}^2 s^2}{(\delta_{ar y})^2}$$

where  $s^2$  and  $\bar{y}$  are estimates from a pilot survey and  $s^2 \approx p(1-p)$  for attributes.

# Example: (blood group)

- 1. What size sample must be drawn from a population of size N = 800in order to estimate the proportion with a given blood group to within 0.04 (i.e. an absolute error of 4%) with probability 0.95?
- 2. What sample size is needed if we know that the blood group is present in no more than 30% of the population?

## Solution:

Read Tutorial 10 Q2.

### THE UNIVERSITY OF SYDNEY STAT3014/3914 Applied Stat.-Sampling C1-Simple random sample

### 1.5 Inference over Subpopulations-Poststratification

**Motivating example:** (dentist) There are 200 children in a village. One dentist takes a simple random sample of 20 and finds 12 children with at least one decayed tooth and a total of 42 decayed teeth. Another dentist quickly checks all 200 children and finds 60 with no decayed teeth.

Estimate the total number of decayed teeth.

 $C_{1} \ge 1 \text{ decayed teeth } C_{2} \text{: no decayed teeth } \text{Total} \\ N_{1} = 140 & N_{2} = 60 & N = 200 \\ n_{1} = 12 & n_{2} = 8 & n = 20 \\ \sum_{i \in C_{1}} y_{i} = 42 = \sum_{i=1}^{n} y'_{i} \\ \end{array}$ 

From a population of N individuals, one simple random sample of n individuals  $y_i$ ,  $i = 1, \dots, n$  is drawn. Separate estimates might be wanted for one of a number of subclasses  $\{C_1, C_2, \dots\}$  which are subsets of the population (sampling frame) using post-stratification.

# Reasons:

- 1. Unavailability of a suitable sampling frame for each stratum even though the stratum sizes  $N_1, \ldots, N_L$  are often obtainable from official statistics.
- 2. Inability to classify population elements into an appropriate stratum without actual contact,

e.g. personal characteristics such as educational level and political preference and household characteristics such as owned/rented accommodation, income level and household size are unknown

- 3. Multi-variate and multi-purpose nature of most surveys.
- 4. Post-stratification is to correct the distorted sample proportion due to non-response.

 $C_1$ 



### Example:

POPULATION (SAMPLING FRAME)SUBPOPULATIONAustralian populationunemployed Queenslandersretailerssupermarketsthe employedthe employed working overtime

**Solution:** The estimate for the *overall* average no. of decay teeth using *overall sample mean* is

The estimate for the overall or conditional total number of decayed teeth,  $\boldsymbol{Y}$  is

ignoring the information of  $n_1 = 12$  and  $N_1 = 140$  from the second dentist. Using these information and condition on *those with at least* one decayed teeth, the average no. of decay teeth is

Alternative estimate for the average no. of decayed teeth using the *conditional sample mean* is

Hence the estimate for the total no. of decayed teeth is

Read Tutorial 10 Q3.



**Formulae:** Denote the total number of items in class  $C_l$  by  $N_l$ . Note that  $N_l$  is generally unknown but we can estimate  $N_l$  by

$$\widehat{N}_l = N n_l / n$$

where  $w_l = n_l/n$  is the sample proportion of units falling into  $C_l$  and the corresponding population proportion is  $W_l = N_l/N$ .

The same technique is used to estimate population mean & total in  $C_l$ :

$$\bar{Y}_l = \sum_{i \in C_l} y_i / N_l$$
, and  $Y_l = \sum_{i \in C_l} y_i$ 

using

Data  $C_1$ 

 $C_2$ 

Mean

Then

 $\left| \bar{y}' = \sum_{i=1}^{n} y'_i / n \right|$  estimates the mean  $\bar{Y}' = \sum_{i=1}^{N} y'_i / N.$ 

The unbiased estimator for the total  $Y_l = N\bar{Y}' = \sum_{i=1}^N Y'_i$  in  $C_l$  and its variance are

$$\widehat{Y}_{pst,lm1} = N\overline{y}' = N\sum_{i=1}^{n} y_i'/n \quad \text{and} \quad \operatorname{var}(\widehat{Y}_{pst,lm1}) = N^2 \left(1 - \frac{n}{N}\right) \frac{s_l'^2}{n}$$

1. When  $N_l$  is known, the unbiased estimator for mean  $\bar{Y}_l = Y_l/N_l$  in  $C_l$  and its variance estimate are

$$\widehat{\bar{Y}}_{pst,lm1} = N\bar{y}'/N_l = \bar{y}'/W_l \quad \text{and} \quad \operatorname{var}(\widehat{\bar{Y}}_{pst,lm1}) = \frac{1}{W_l^2} \left(1 - \frac{n}{N}\right) \frac{{s_l'}^2}{n}$$



where  $S'_l^2$  can be estimated by  $s'_l^2$ :

$$S_l'^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i' - \bar{Y}')^2$$
 and  $s_l'^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i' - \bar{y}')^2$ .

**Note:** the random variable  $n_l$  are not included in these calculations.

2. When  $N_l$  is unknown, we can estimate the total  $Y_l$  by  $N\bar{y}'$  but we cannot estimate  $\bar{Y}_l$  by  $\bar{y}'/W_l$ . A natural way is to estimate  $W_l = \frac{N_l}{N}$  by  $w_l = \frac{n_l}{n}$ :

$$\widehat{\bar{Y}}_{pst,lm2} = n\bar{y}'/n_l = \bar{y}_l$$
 and  $\operatorname{var}(\widehat{\bar{Y}}_{pst,lm2}) \simeq \frac{1}{W_l} \left(1 - \frac{n}{N}\right) \frac{{s_l}^2}{n}$ 

Similarly the total estimator and its variance estimate are

$$\begin{aligned} \widehat{Y}_{pst,lm2} &= N_l \bar{y}_l = N_l n \bar{y}'/n_l \quad \text{and} \quad \operatorname{var}(\widehat{Y}_{pst,lm2}) \simeq N^2 W_l \left(1 - \frac{n}{N}\right) \frac{s_l^2}{n} \\ \text{where } s_l^2 &= \frac{1}{n_l - 1} \sum_{i \in C_l} (y_i - \bar{y}_l)^2 \text{ is the sample variance in } C_l \text{ and} \\ \frac{N_l^2}{W_i} &= N^2 W_l. \end{aligned}$$

**Theorem:** For the estimator  $\hat{\bar{Y}}_l = \bar{y}_l = (n/n_l)\bar{y}'$ ,

$$E(\bar{y}_l) = \bar{Y}_l$$
  

$$Var(\bar{y}_l) \simeq \frac{1}{W_l} \left(1 - \frac{n}{N}\right) \frac{S_l^2}{n}$$

provided we define  $E(\bar{y}_l) = \bar{Y}_l$  when  $n_l = 0$ .



**Proof:** Using conditional expectation,

$$E(\bar{y}_l) = E_{n_l} \{ E(\bar{y}_l | n_l) \}$$
  
=  $E(\bar{y}_l | n_l = 0) \operatorname{Pr}(n_l = 0) + \sum_{i \ge 1} E(\bar{y}_l | n_l = i) \operatorname{Pr}(n_l = i)$   
=  $\bar{Y}_l \operatorname{Pr}(n_l = 0) + \sum_{i \ge 1} \bar{Y}_l \operatorname{Pr}(n_l = i) = \bar{Y}_l$ 

where for each fixed  $n_l$ , a simple random sample of size  $n_l$  is drawn from  $C_l$ . Further using conditional variance,

$$\operatorname{Var}(\bar{y}_l) = \operatorname{Var}(\underbrace{E(\bar{y}_l|n_l)}_{\bar{Y}_l} + E(\operatorname{Var}(\bar{y}_l|n_l)) = 0 + E(\operatorname{Var}(\bar{y}_l|n_l))$$
$$= E\left[\left(1 - \frac{n_l}{N_l}\right)\frac{S_l^2}{n_l}\right] = E\left(\frac{S_l^2}{n_l} - \frac{S_l^2}{N_l}\right) \simeq \frac{S_l^2}{nW_l} - \frac{S_l^2}{NW_l}$$
$$= \frac{1}{W_l}\left(1 - \frac{n}{N}\right)\frac{S_l^2}{n}$$

since

 $E\left(\frac{1}{n_l}\right) \simeq \frac{1}{nW_l} + \frac{1-W_l}{n^2W_l^2} \simeq \frac{1}{nW_l}$  where the ignored term: TT7

extra var. = 
$$\frac{1 - W_l}{n^2 W_l^2} S_l^2$$

is the extra variability due to the random sample size  $n_l$ .

Theorem: 
$$E\left(\frac{1}{n_l}\right) \simeq \frac{1}{nW_l} + \frac{1 - W_l}{n^2 W_l^2}.$$
  
Proof:  $\frac{1}{n_l} = (n_l - nW_l + nW_l)^{-1} = \frac{1}{nW_l} \left(1 + \frac{n_l - nW_l}{nW_l}\right)^{-1}$   
 $= \frac{1}{nW_l} \left[1 - \left(\frac{n_l - nW_l}{nW_l}\right) + \left(\frac{n_l - nW_l}{nW_l}\right)^2 - \dots\right]$ 



Hence 
$$E\left(\frac{1}{n_l}\right) = \frac{1}{nW_l} \left[1 - \frac{E(n_l - nW_l)}{nW_l} + \frac{E(n_l - nW_l)^2}{n^2W_l^2} - \dots\right]$$
  
 $\simeq \frac{1}{nW_l} + \frac{1}{nW_l} \left(\frac{nW_l(1 - W_l)}{n^2W_l^2}\right) = \frac{1}{nW_l} + \frac{1 - W_l}{n^2W_l^2}$ 

since  $E(n_l) = nW_l$  and  $Var(n_l) = nW_l(1 - W_l)$ .

**Corollary:** The estimator  $\widehat{Y}_{pst,lm2} = N_l \, \overline{y}_l$  is unbiased for  $Y_l$  and has  $\operatorname{Var}(\widehat{Y}_l) = N_l^2 \operatorname{Var}(\overline{y}_l)$ .

When  $N_l$  is known, the estimator  $\widehat{Y}_{pst,lm2} = N_l \, \overline{y}_l$  uses more information (both  $N_l \& n_l$ ) than  $N \overline{y}'$  and so is a better unbiased estimator.

Using this method, the *overall* total and mean estimates are:

$$\widehat{Y}_{pst} = N \sum_{l=1}^{L} W_l \, \overline{y}_l \text{ and } \operatorname{var}\left(\widehat{Y}_{pst}\right) \simeq N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} \sum_{l=1}^{L} W_l s_l^2$$

and

$$\widehat{\bar{Y}}_{pst} = \sum_{l=1}^{L} W_l \, \bar{y}_l \text{ and } \operatorname{var}\left(\widehat{\bar{Y}}_{pst}\right) \simeq \left(1 - \frac{n}{N}\right) \, \frac{1}{n} \sum_{l=1}^{L} W_l s_l^2$$

respectively since

$$\begin{split} \widehat{\bar{Y}}_{pst} &= \frac{1}{N} \sum_{l=1}^{L} \widehat{Y}_{pst,lm2} = \frac{1}{N} \sum_{l=1}^{L} N_l \, \bar{y}_l = \sum_{l=1}^{L} W_l \, \bar{y}_l \\ \operatorname{var}\left(\widehat{\bar{Y}}_{pst}\right) &= \sum_{l=1}^{L} W_l^2 \operatorname{var}(\bar{y}_l) \simeq \sum_{l=1}^{L} W_l^2 \frac{1}{W_l} \left(1 - \frac{n}{N}\right) \frac{s_l^2}{n} \simeq \left(1 - \frac{n}{N}\right) \frac{1}{n} \sum_{l=1}^{L} W_l s_l^2 \\ \operatorname{Compare} \, \widehat{\bar{Y}}_{pst} \text{ to } \quad \widehat{\bar{Y}} = \bar{y} = \sum_{l=1}^{L} w_l \, \bar{y}_l \text{ and } \operatorname{var}\left(\widehat{\bar{Y}}\right) = \left(1 - \frac{n}{N}\right) \frac{s^2}{n} \\ \text{for one SRS without post-stratification, we have} \end{split}$$



 $W_l$  replaces  $w_l$  and  $\sum_{l=1}^{L} W_l s_l^2$  replaces  $s^2$  to correct sample proportions.

Read Tutorial 11 Q1,c,d.

### Post-stratified estimator based on 1 SRS

Par.	$N_l$	Estimator	Variance
$\bar{Y}_l$	known	$\widehat{\bar{Y}}_{pst,lm1} = N\bar{y}'/N_l$	$\operatorname{var}(\widehat{\bar{Y}}_{pst,lm1}) = \frac{1}{W_l^2} \left(1 - \frac{n}{N}\right) \frac{{s'_l}^2}{n}$
$Y_l$	unknown	$\widehat{Y}_{pst,lm1} = N\bar{y}'$	$\operatorname{var}(\widehat{Y}_{pst,lm1}) = N^2 \left(1 - \frac{n}{N}\right) \frac{{s'_l}^2}{n}$
$ar{Y_l}$	unknown	$\hat{\bar{Y}}_{pst,lm2} = \bar{y}_l$	$\operatorname{var}(\widehat{\bar{Y}}_{pst,lm2}) \simeq \frac{1}{W_l} \left(1 - \frac{n}{N}\right) \frac{{s_l}^2}{n}$
$Y_l$	known	$\widehat{Y}_{pst,lm2} = N_l \bar{y}_l$	$\operatorname{var}(\widehat{Y}_{pst,lm2}) \simeq N^2 W_l \left(1 - \frac{n}{N}\right) \frac{{s_l}^2}{n}$
$\overline{Y}$	known	$\widehat{\overline{Y}}_{pst} = \sum_{l=1}^{L} W_l \overline{y}_l$	$\operatorname{var}(\widehat{\overline{Y}}_{pst}) \simeq \left(1 - \frac{n}{N}\right) \sum_{l=1}^{L} W_l \frac{s_l^2}{n}$
Y	known	$\widehat{Y}_{pst} = \sum_{l=1}^{L} N_l \overline{y}_l$	$\operatorname{var}(\widehat{Y}_{pst}) \simeq N^2 \left(1 - \frac{n}{N}\right) \sum_{l=1}^{L} W_l \frac{s_l^2}{n}$

where  $y'_i = y_i$  if  $i \in C_l \& 0$  otherwise,

$$\bar{y}' = \frac{1}{n} \sum_{i \in C_l} y_i, \ \bar{y}_l = \frac{1}{n_l} \sum_{i \in C_l} y_i,$$

$$s_l'^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i' - \bar{y}')^2$$
, and  $s_l^2 = \frac{1}{n_l - 1} \sum_{i \in C_l} (y_i - \bar{y}_l)^2$ .