

Statistical modeling of data collected from a stratified random sampling.

- Data: $(\underline{X}_{1,1}, \underline{X}_{2,1}, \dots, \underline{X}_{n_1,1}), \dots, (\underline{X}_{1,L}, \underline{X}_{2,L}, \dots, \underline{X}_{n_L,L})$, $\in \underline{\mathbb{S}_1} \in \underline{\mathbb{S}_L}$

sampling probability in \mathbb{S}_l : $1/N_l$

depends on

- N_l 's
- $F_{0,l}$'s

where $(\underline{X}_{1,l}, \dots, \underline{X}_{n_l,l})$, $l = 1, \dots, L$, is the data collected from the s.r.s. (either with or without replacement) taken within the l th stratum $\underline{\mathbb{S}_l}$.

distribution of data

subpopulation

check
Def. 19
(LNp.62)

- $(\underline{X}_{1,l}, \dots, \underline{X}_{n_l,l})$: since a s.r.s. is taken within each stratum, the joint distribution of the data from the stratum $\underline{\mathbb{S}_l}$ is as that given in LNp.11-12, with F_0 replaced by $F_{0,l}$

Q: What is the joint dist. of all data?

data from different strata are independent

$$\prod_{l=1}^L P(\underline{x}_{1,2}, \dots, \underline{x}_{n_l,2})$$

Definition 20 (some intuitive estimators of the parameters of population and stratum)

subpopulation $\underline{\mathbb{S}_l}$: since a s.r.s. is taken within each stratum,

check
Note8
(LNp.28)

– mean $\underline{\mu}_l$: estimated by the subsample mean $\underline{\bar{X}}_l \equiv \frac{1}{n_l} \sum_{k=1}^{n_l} \underline{X}_{k,l}$

– total $\underline{\tau}_l$: estimated by subsample total $\underline{T}_l \equiv N_l \underline{\bar{X}}_l$

(unbiased)
sample
variance

– variance $\underline{\sigma}_l^2$: estimated by $\underline{s}_l^2 \equiv \frac{1}{n_l - 1} \sum_{k=1}^{n_l} (\underline{X}_{k,l} - \underline{\bar{X}}_l)^2$ under with replacement, and by $\left(1 - \frac{1}{N_l}\right) \underline{s}_l^2$ under without replacement

- (whole) population Ω : under a stratified random sample,

– mean $\underline{\mu}$: estimated by the stratified sample mean

$$\underline{\bar{X}}_{\mathbb{S}} \equiv \frac{1}{N} \sum_{l=1}^L N_l \underline{\bar{X}}_l = \sum_{l=1}^L W_l \underline{\bar{X}}_l = \frac{1}{N} \sum_{l=1}^L \frac{1}{(n_l/N_l)} \left(\sum_{k=1}^{n_l} \underline{X}_{k,l} \right),$$

Thm20
(LNp.61)

Why should
not equal?
check the
graph in
LNp.60

since $\underline{\mu} = \sum_{l=1}^L W_l \underline{\mu}_l$. τ_l stratum $\frac{N_l}{N}$ fraction cf. weights

Note. $\underline{\bar{X}}_{\mathbb{S}} \neq \frac{1}{n} \sum_{l=1}^L \sum_{k=1}^{n_l} \underline{X}_{k,l} = \sum_{l=1}^L \frac{n_l}{n} \underline{\bar{X}}_l$ in general,

weights they are equal only when $\frac{n_l}{n} = \frac{N_l}{N}$, $l = 1, \dots, L$.)

– total τ ($= N \underline{\mu}$): estimated by $\underline{T}_{\mathbb{S}} \equiv N \underline{\bar{X}}_{\mathbb{S}}$ iff $\frac{N}{N} = \frac{n_l}{N_l}$ sampling fraction ↑ weight ↓

This explains
why in LNp.63
within-stratum
variations
should be
small

– FYI. An intuitive estimator of the population variance σ^2 can be developed, based on the relation between σ^2 and $\underline{\mu}_l$'s, $\underline{\sigma}_l^2$'s (Thm. 20, LNp.61), by using the estimators $\underline{\bar{X}}_l$'s and \underline{s}_l^2 's (or $(1 - \frac{1}{N_l}) \underline{s}_l^2$'s).

Theorem 22 (mean and variance of the stratified estimator of population mean)

- Under stratified random sampling, with or without replacement, $E(\underline{\bar{X}}_{\mathbb{S}}) = \underline{\mu}$.
- Under stratified random sampling,

– with replacement, $Var(\underline{\bar{X}}_{\mathbb{S}}) = \sum_{l=1}^L W_l^2 \left(\frac{\underline{\sigma}_l^2}{n_l} \right)$. stratum $\frac{N_l}{N}$ fraction unbiased

– without replacement, $Var(\underline{\bar{X}}_{\mathbb{S}}) = \sum_{l=1}^L W_l^2 \left(\frac{\underline{\sigma}_l^2}{n_l} \right) \left(1 - \frac{n_l - 1}{N_l - 1} \right)$. Var($\underline{\bar{X}}_l$) Note8, LNp.28

Proof: The expectation of the stratified estimator \bar{X}_S is

$$E(\bar{X}_S) = E\left(\sum_{l=1}^L W_l \bar{X}_l\right) = \sum_{l=1}^L W_l E(\bar{X}_l) = \sum_{l=1}^L W_l \mu_l = \mu.$$

Thm 20, LNp.61

Ch7, p.66

indep.

Since the data from different strata are independent of one another, the subsample means $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_L$ are independent random variables, and

$$\begin{aligned} S_1 &\rightarrow (x_{1,1}, x_{2,1}, \dots, x_{n_1,1}) \\ S_2 &\rightarrow (x_{1,2}, x_{2,2}, \dots, x_{n_2,2}) \\ S_L &\rightarrow (x_{1,L}, x_{2,L}, \dots, x_{n_L,L}) \end{aligned} \quad \text{Var}(\bar{X}_S) = \text{Var}\left(\sum_{l=1}^L W_l \bar{X}_l\right) = \sum_{l=1}^L W_l^2 \text{Var}(\bar{X}_l).$$

∴ S.R.S. in each strata

Since S.R.S. is taken within each stratum, the results follows from Thm.2 (LNp.17) and Thm.3 (LNp.18) respectively for with and without replacement.

Note 17 (Some notes about the mean and variance of the stratified estimator of μ)

- Under stratified random sampling, \bar{X}_S is an unbiased estimator of μ .
- If the sampling fractions (i.e., (n_l/N_l) 's) within all strata are small, then with replacement \approx without replacement, and

similar probabilities

$$n_l \ll N_l, l=1, \dots, N \quad \sum_{l=1}^L W_l^2 \left(\frac{\sigma_l^2}{n_l}\right) \approx \sum_{l=1}^L W_l^2 \left(\frac{\sigma_l^2}{n_l}\right) \left(1 - \frac{n_l-1}{N_l-1}\right) \approx 1$$

Definition 22 (estimated standard error of the stratified estimator of population mean)

- Under stratified random sampling with replacement, since s_l^2 is an unbiased estimator of σ_l^2 , the $\text{Var}(\bar{X}_S)$ can be estimated by

$$\text{unbiased} \rightarrow s_{\bar{X}_S}^2 = \sum_{l=1}^L W_l^2 \left(\frac{s_l^2}{n_l}\right). \quad \text{Def. 20, LNp. 64} \quad \therefore \text{s.r.s in each strata and Thm 9 (LNp.24)}$$

Def. 20, LNp. 64

∴ s.r.s in each strata and Thm 9 (LNp.24)

Def. 20, LNp. 64

Stratum	N_l	$W_l = \frac{N_l}{N}$	μ_l	increasing	σ_l	$\leq \sigma$
A	98	0.249	182.9		103.4	smallest
B	98	about equal	526.5		204.8	
C	98	about equal	956.3		243.5	
D	99	0.252	1591.2		419.2	largest

x: unknown
in sampling
survey

good strata?
Check "Recall"
in LNp.63

- For a without-replacement stratified random sample of size n , suppose we choose $n_1 = n_2 = n_3 = n_4 = n/4$. Neglecting the finite population correction, we have

cf.

Thm22, LNp.65

1/n_g

$$\sigma_{\bar{X}_S} = \sqrt{\frac{4}{n} \sum_{l=1}^4 W_l^2 \sigma_l^2} = \frac{268.4}{\sqrt{n}}$$

$$\begin{aligned} ① & \because n_g \ll N_g, n \ll N \\ ② & \because 1 - \frac{n_g - 1}{N_g - 1} \approx 1 - \frac{n - 1}{N - 1} \end{aligned}$$

- For a without-replacement s.r.s. of size n , neglecting the finite population correction, we have (see Ex.4, LNp.20)

population variance $\sigma^2 = 589.7^2$

Thm3, LNp.18

$$\sigma_{\bar{X}} = \frac{589.7}{\sqrt{n}}$$

- Note that the stratification has resulted in a tremendous gain in precision: $\sigma_{\bar{X}_S} \approx 0.455 \times \sigma_{\bar{X}}$ $\Rightarrow \sigma_{\bar{X}_S}^2 / \sigma_{\bar{X}}^2 = 0.207$. The stratified estimator \bar{X}_S based on a total sample size of $n/5$ is as precise as \bar{X} based on a s.r.s. of size n . (cf. the reduction in variance due to stratification is comparable to that achieved by using a ratio estimator given in Ex.18, LNp.58). LNp.57

• Methods of allocation in stratified random sampling

- Q: Why and when can a stratification produce a dramatic improvement in precision?

• Why? :: exclude many biased samples

• When? i.e., n_g's = ? strata = ?