

Stratified Sampling

Idea: divide heterogeneous

population to strata,

which are homogeneous within each

strata and heterogeneous between

strata

1) $N = \sum_{i=1}^k N_i$ divide into k strata

2) n_i sample size in each strata
(SRS WOR)

$$n = \sum_{i=1}^k n_i$$

⊛ Cluster sampling:

within heterogeneous

between homogeneous

Stratified (variance reduction technique)

$E(Y)$

$$n : \frac{1}{n} \sum y_i = \bar{y} \quad \bar{y} - \text{unbiased}$$

$$\text{var}(\bar{y}) = \frac{\text{var}(y)}{n}$$

K subgroups : w_k - prob. to be in group k

$$n_k = w_k \cdot n$$

$$\hat{y} = \sum w_k \cdot \bar{y}_k$$

\bar{y}_k - s. mean within k^{th} group

$$E(\bar{y}_{\text{strat}}) = \sum_{k=1}^K w_k E(\bar{y}_k) = \sum_{k=1}^K w_k \mu_k = \mu$$

$$\text{var}(\bar{y}_{\text{strat}}) = \sum_{k=1}^K w_k^2 \cdot \text{var}(\bar{y}_k) =$$

$$= \sum_{k=1}^K \frac{w_k^2}{n^2} \cdot \frac{1}{\cancel{w_k}} \cdot \sigma_k^2 = \frac{1}{n} \sum_{k=1}^K w_k \cdot \sigma_k^2$$

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

$$\text{Var}(\bar{y}) = \sigma^2/n$$

$$\text{Var}(\bar{y}) \stackrel{?}{=} \text{Var}(\hat{y}_{\text{strat}})$$

z - strat.
covariate

$$\begin{aligned} \text{Var}(y) &= E(\text{Var}(y|z)) + \text{Var}(E(y|z)) = \\ &= E\left(\sum \sigma_k^2 \cdot I(z=k)\right) + \text{Var}\left(\sum \mu \cdot I(z=k)\right) = \\ &= \sum \sigma_k^2 \cdot E(I(z=k)) + E\left(\left(\sum \mu \cdot I(z=k)\right)^2\right) - \\ &\quad E^2\left(\sum \mu \cdot I(z=k)\right) = \\ &= \sum \sigma_k^2 \cdot w_k + \sum \mu_k^2 w_k - \mu^2 \\ &= \sum \sigma_k^2 w_k + \sum w_k (\mu_k - \mu)^2 \end{aligned}$$

within

between

$$\begin{aligned} \text{var}(\bar{y}) &= \frac{1}{n} \sum \sigma_k^2 w_k + \frac{1}{n} \sum w_k (\mu_k - \mu)^2 \\ &\quad \text{"} \\ &\quad \text{Var}(\hat{y}_{\text{strat}}) \end{aligned}$$

How to chose n_i , $i=1, \dots, K$

(i) min cost of survey

(ii) max precision

1) Equal allocation:

$$n_i = \frac{n}{K}$$

2) Proportional allocation:

$$n_i \propto N_i$$

$$n_i = \delta N_i$$

$$n = \delta \cdot N \Rightarrow \delta = n/N$$

$$n_i = \left(\frac{n}{N}\right) \cdot N_i$$

$$\text{Sampling error } (\hat{y}_{\text{prop}}) \leq \text{s.e. } (\hat{y}_{\text{SRS}})$$

more diverse \Rightarrow more precise $\hat{y}_{\text{prop.}}$

3) Optimal (Neyman) allocation

$$h_i \propto N_i s_i$$

$$h_i = \delta^* N_i s_i$$

$$n = \delta^* \sum N_i s_i \Rightarrow \delta^* = \frac{n}{\sum N_i s_i}$$

$$h_i = \frac{n}{\sum N_i s_i} N_i s_i$$

\hat{y}_{opt} As precise / more precise than $\hat{y}_{prop.}$
more precise when stratas' se. differ

Cost of survey:

$$C = C_0 + \sum C_i h_i$$

\uparrow per unit in i th state

$$\text{Var}(\bar{y}) \rightarrow \min_{n_i} \quad \text{s.t. } C = C_0 + \sum C_i h_i$$

$$L = \underset{\text{stat}}{\text{Var}(\bar{y})} + \lambda^2 (C - C_0) =$$

$$= \sum w_i \left(\frac{1}{h_i} - \frac{1}{N_i} \right) S_i^2 + \lambda^2 \cdot \sum C_i h_i$$

$$= \sum \frac{w_i S_i^2}{n_i} + \lambda^2 \sum C_i h_i - \sum \frac{w_i S_i^2}{N_i} =$$

$$= \sum \left(\frac{w_i S_i}{\sqrt{n_i}} - \lambda \sqrt{C_i n_i} \right)^2 + B$$

$$L \rightarrow \min \quad \text{at} \quad \frac{w_i S_i}{\sqrt{n_i}} = \lambda \sqrt{C_i n_i}$$

$$h_i = \frac{1}{\lambda} \cdot \frac{w_i S_i}{\sqrt{C_i}}$$

λ \nearrow min var, fixed cost (i)
 \searrow min cost, fixed var (ii)

(i) $C = C^*$

$$\sum C_i h_i = C^*$$

$$\sum C_i \cdot \frac{w_i S_i}{\lambda \sqrt{C_i}} = C^*$$

$$\lambda_{FC} = \frac{\sum \sqrt{C_i} \cdot w_i S_i}{C^*}$$

$$h_i = \frac{1}{\lambda_{FC}} \cdot \frac{w_i S_i}{\sqrt{C_i}}$$

(ii) $\sum \left| \frac{1}{h_i} - \frac{1}{N_i} \right| w_i^2 S_i^2 = V_0$

$$\sum \frac{\lambda \sqrt{C_i}}{w_i S_i} w_i^2 S_i^2 = V_0 + \sum \frac{w_i^2 J_i^2}{N_i}$$

$$\lambda_{FV} = \frac{V_0 + \sum \frac{w_i^2 S_i^2}{N_i}}{\sum w_i S_i \sqrt{C_i}}$$

$$h_i = \frac{1}{\lambda_{FV}} \cdot \frac{w_i S_i}{\sqrt{C_i}}$$

sample size
SRS

$$\frac{\bar{X}}{e^2} = \frac{z_{\alpha/2}^2 \cdot \sigma^2}{e^2}$$

$$\frac{\hat{p}}{e^2} = \frac{z_{\alpha/2}^2 \cdot p \cdot (1-p)}{e^2}$$

prop

$$\frac{z_{\alpha/2}^2 \cdot \sum W_k \cdot \sigma_k^2}{e^2}$$

$$\frac{z_{\alpha/2}^2 \cdot \sum W_k \cdot p_k (1-p_k)}{e^2}$$

opt

$$\frac{z_{\alpha/2}^2 \cdot (\sum W_k \cdot \sigma_k)^2}{e^2}$$

$$\frac{z_{\alpha/2}^2 \cdot \left(\sum W_k \sqrt{p_k \cdot (1-p_k)} \right)^2}{e^2}$$

if $\sigma_k = \sigma_j$, $n_k = n_j$ $\forall k, j$

$$\frac{(\sum W_k \sigma_k)^2}{e^2} = \frac{\left(\sum \frac{1}{k} \cdot \sigma \right)^2}{e^2} = \frac{k/k \cdot \sigma^2}{e^2}$$