

Math 152 2/5/09.

Chptr 7 #49

$N = 1000$ (Quadrats) $\left(\begin{array}{l} \text{Square all finite} \\ \text{population corrections} \\ \text{and always treat } \frac{1}{n-1} \text{ as } \frac{1}{n} \end{array} \right)$
 $n = 50$

$X_i =$ area covered by vegetation in i^{th} quadrat

$Y_i =$ # of birds in i^{th} quadrat.

$$\sum X_i = 3000$$

$$\bar{X} = \frac{3000}{50} = 60$$

$$\sum Y_i = 150$$

$$\bar{Y} = \frac{150}{50} = 3$$

$$\sum X_i^2 = 225,000$$

$$\hat{\sigma}_X^2 = \frac{225000}{50} - (60)^2 = 900$$

$$\sum Y_i^2 = 650$$

$$\hat{\sigma}_Y^2 = \frac{650}{50} - 3^2 = 4$$

$$\sum X_i Y_i = 11,000$$

$$\hat{\sigma}_{XY} = \frac{11000}{50} - (60)(3) = 40.$$

a) $R = \frac{\bar{Y}}{\bar{X}} = \frac{3}{60} = .05.$

b) $S_R^2 \approx \frac{1}{50} \frac{1}{(60)^2} \left((.05)^2 \cdot 900 + 4 - 2(.05)(40) \right)$

$$S_R \approx .00354.$$

b) continued

$$\text{Sample variance of } \bar{X} \approx \frac{\hat{\sigma}_x^2}{n}$$

$$\text{Standard error of } \bar{X} \approx \frac{\hat{\sigma}_x}{\sqrt{n}} = \frac{30}{\sqrt{50}} \approx 4.24.$$

90% confidence interval for μ_x

$$\approx 60 \pm (1.64)(4.25) \approx (53.05, 66.95)$$

$$\text{Sample variance of } \bar{Y} \approx \frac{\hat{\sigma}_y^2}{n}$$

$$\text{Standard error of } \bar{Y} \approx \frac{\hat{\sigma}_y}{\sqrt{n}} = \frac{2}{\sqrt{50}} \approx .283$$

90% confidence interval for
 μ_y

$$\approx 3 \pm (1.64)(.283) \approx (2.53, 3.46)$$

$$c) T = N\bar{Y} = 3000$$

$$\text{Var}(T) = N^2 \text{Var}(\bar{Y})$$

$$\text{Standard error of } T = N \cdot \text{standard error}(\bar{Y}) \\ \approx 1000(.283) = 283$$

a 95% confidence interval for T
 is
 $3000 \pm (1.96)(283)$
 $= (2445.32, 3554.68)$

d). If the total area covered by vegetation τ_x is known then $\mu_x = \frac{\tau_x}{1000}$ so

$$\bar{Y}_R = \mu_x \frac{\bar{Y}}{\bar{X}} \text{ is known}$$

and $T_R = N \bar{Y}_R$ gives an estimate
 of τ_y .

Now

$$\text{Var}(T_R) = N^2 \text{Var}(\bar{Y}_R)$$

$$\approx N^2 \cdot \frac{1}{n} (r^2 \sigma_x^2 + \sigma_y^2 - 2r \sigma_{xy})$$

$$\approx \frac{1000^2}{50} ((.05)^2 \cdot 900 + 4 - 2(.05)(.40)) = 4500$$

Now

$$\text{but } \text{Var}(T) \approx (283)^2 \approx 80,000.$$

so if the bias is small, then T_R will be a better estimate.

$$E(\bar{Y}_R) - \mu_Y \approx \frac{1}{50} \cdot \frac{1}{\mu_X} ((.05)900 - 40)$$

and from part b we have

"90% confidence" that

$$\frac{1}{50} \cdot \frac{1}{\mu_X} ((.05)(900) - 40) \leq$$

$$\frac{1}{50} \cdot \frac{1}{53} ((.05)(900) - 40) \approx .002$$

and therefore that the bias in T_R

is roughly 2.

We expect the mean squared error of

T_R to be roughly $45000 + 4$

while the mean squared error of T is
 $\approx 80,000.$

Stratified Sampling.

$$N = N_1 + \dots + N_L$$

$$\mu_l = \frac{1}{N_l} \sum_{i=1}^{N_l} X_{il}$$

$$\mu = \sum_{l=1}^L W_l \mu_l, \quad W_l = \frac{N_l}{N}.$$

$$\bar{X}_l = \frac{1}{n_l} \sum_{i=1}^{n_l} X_{il} \quad l=1, \dots, L.$$

$$n_1 + \dots + n_L = n.$$

$$\bar{X}_s = \sum_{l=1}^L W_l \bar{X}_l \quad \left(\begin{array}{l} \text{stratified} \\ \text{sample} \\ \text{mean} \end{array} \right)$$

We have:

$$E(\bar{X}_s) = \mu$$

$$\begin{aligned} \text{and } \text{Var}(\bar{X}_s) &= \sum_{l=1}^L W_l^2 \text{Var}(\bar{X}_l) \\ &= \sum_{l=1}^L W_l^2 \frac{1}{n_l} \left(1 - \frac{n_l-1}{N_l-1}\right) \sigma_l^2. \end{aligned}$$

and

$$\text{Var}(\bar{X}_s) \approx \sum_{l=1}^L \frac{W_l^2 \sigma_l^2}{n_l}$$

neglecting finite population corrections.

e.g. Stratify the population of hospitals according to # of beds.

stratum	N_l	W_l	μ_l	σ_l
1-98	98	.249	182.9	103.4
99-196	98	.249	526.5	204.8
197-294	98	.249	956.3	243.5
295-393	99	.251	1591.2	419.2

n = sample size.

$$n_1 = n_2 = n_3 = n_4 = n/4.$$

$$\begin{aligned} \text{Var}(\bar{X}_s) &\approx \sum_{l=1}^4 \frac{W_l^2 \sigma_l^2}{n/4} = \frac{4}{n} \sum_{l=1}^4 W_l^2 \sigma_l^2 \\ &= \frac{72,042.6}{n} \end{aligned}$$

$$\sigma_{\bar{X}_s} = \frac{268.4}{\sqrt{n}}$$

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

$$\frac{268.4}{\sqrt{n_1}} = \frac{589.7}{\sqrt{n_2}}$$

$$\Leftrightarrow n_1 = \left(\frac{268.4}{589.7} \right)^2 n_2 \quad \hat{n} \approx .2 n_2.$$

For population totals we have:

$$\bar{T}_s = N \bar{X}_s \quad E(T_s) = \tau$$

$$\text{Var}(T_s) = N^2 \text{Var}(\bar{X}_s)$$

$$= \sum_{e=1}^L W_e^2 \left(\frac{1}{n_e} \right) \left(1 - \frac{n_e-1}{N_e-1} \right) \sigma_e^2$$

To get practically useful expressions we must estimate the σ_e^2 by

$$s_e^2 = \frac{1}{n_e-1} \sum_{i=1}^{n_e} (X_{ie} - \bar{X}_e)^2$$

and then

$\text{Var}(\bar{X}_s)$ is estimated as:

$$S_{\bar{X}_s}^2 = \sum_{e=1}^L W_e^2 \left(\frac{1}{n_e} \right) \left(1 - \frac{n_e}{N_e} \right) s_e^2$$

Allocation for stratified sampling

"Optimal" allocation

Minimizing

$$\text{Var}(\bar{X}_s) = \sum_{l=1}^L \frac{w_l^2 \sigma_l^2}{n_l}$$

(as a function of n_1, \dots, n_L)

subject to the constraint

$$\text{that } n_1 + \dots + n_L = 1.$$

Lagrange Multipliers:

$$\left[\nabla (\text{Var}(\bar{X}_s)) \right]_l = - \frac{w_l^2 \sigma_l^2}{n_l^2}$$

$$\left[\nabla (n_1 + \dots + n_L) \right]_l = 1.$$

Solve the system.

$$- \frac{w_l^2 \sigma_l^2}{n_l^2} = 1 \cdot \lambda \quad l=1, \dots, L$$

$$n_1 + \dots + n_L = n.$$

(we may as well use $-\lambda$).
so solve.

$$W_l \sigma_l^2 = n_l^2 \lambda. \quad l = 1, \dots, L.$$

(+ constraint).

$$n_l = \frac{W_l \sigma_l}{\sqrt{\lambda}}$$

$$n = \sum_{l=1}^L \frac{W_l \sigma_l}{\sqrt{\lambda}}$$

$$\frac{1}{\sqrt{\lambda}} = \frac{n}{\sum_{l=1}^L W_l \sigma_l}$$

$$n_l = \frac{W_l \sigma_l \cdot n}{\sum_{l=1}^L W_l \sigma_l}$$

With this choice for n_l .

we call the estimator

\bar{X}_{50} (stratified optimal).

and

$$\begin{aligned} \text{Var}(\bar{X}_{50}) &\approx \sum_{l=1}^L \frac{W_l^2 \sigma_l^2}{\frac{W_l \sigma_l \cdot n}{\sum_{l=1}^L W_l \sigma_l}} \\ &= \frac{1}{n} \left(\sum_{l=1}^L W_l \sigma_l \right)^2. \end{aligned}$$

e.g. Same strata for hospitals as before

	A	B	C	D
	1-98	99-196	197-294	295-393
$\frac{W_l \sigma_l}{\sum W_l \sigma_l}$.106	.210	.250	.434

Optimal allocation is difficult to achieve in practice and since it depends on variability, it depends on what question is being asked.

So if your survey asks more than one

A practical alternative is:

Proportional allocation

$$\frac{n_1}{N_1} = \frac{n_2}{N_2} = \dots = \frac{n_L}{N_L}$$

i.e. $n_e = n W_e$.

$$\bar{X}_{sp} = \sum_{e=1}^L W_e \bar{X}_e$$

is the unweighted mean of all sample values.

We have:

$$\text{Var}(\bar{X}_{sp}) \approx \frac{1}{n} \sum_{e=1}^L W_e^2 \sigma_e^2$$

(ignoring finite pop.).

$$= \sum_{e=1}^L W_e^2 \text{Var}(\bar{X}_e) \approx \sum_{e=1}^L W_e^2 \frac{\sigma_e^2}{n_e} = \frac{1}{n} \sum_{e=1}^L \frac{W_e \sigma_e^2}{W_e}$$

With some algebra, we have

$$\text{Var}(\bar{X}_{sp}) - \text{Var}(\bar{X}_{so})$$

$$= \frac{1}{n} \sum_{l=1}^L W_l (\sigma_l - \bar{\sigma})^2$$

$$\text{where } \bar{\sigma} = \sum_{l=1}^L W_l \sigma_l.$$

Ans If all σ_l are equal the two methods give the same result. If the σ_l are varying widely between strata then the optimal allocation is "better".

We can also compare S.R.S. with proportional as follows:

$$\text{Var}(\bar{x}) - \text{Var}(\bar{x}_{sp}) = \frac{1}{n} \sum_{l=1}^L W_l (\mu_l - \mu)^2$$

(after some algebra).

so proportional sampling is to be preferred when the means of the strata are variable.

For our hospitals example we have.

$$\text{Var}(\bar{x}_{sp}) = \text{Var}(\bar{x}_{so}) + \frac{1}{n} \sum W_l (\sigma_l - \bar{\sigma})^2$$

so that

$$\frac{\text{Var}(\bar{x}_{sp})}{\text{Var}(\bar{x}_{so})} = 1 + \frac{\sum W_l (\sigma_l - \bar{\sigma})^2}{(\sum W_l \sigma_l)^2}$$

$$= 1 + 0.218.$$

and

$$\frac{\text{Var}(\bar{X}_{\text{srs}})}{\text{Var}(\bar{X}_{\text{sp}})} = 1 + \frac{\sum w_e (\mu_e - \bar{\mu})^2}{\sum w_e \sigma_e^2}$$
$$= 1 + 3.83.$$

This situation is typical in applications.

The gain in using proportional allocation rather than S.P.S. is large compared to the gain in using optimal allocation over proportional allocation.