

- an argument based on the CLT can be used to show that \underline{R} is approximately normally distributed, i.e., $\underline{R} \stackrel{D}{\approx} \underline{N}(\mu_R, \sigma_R^2)$, when sample size n is large.
- Applications** ($\because \text{bias} \rightarrow 0$) $\underline{r}_{xy} \xrightarrow{n \rightarrow 0}$ can be replaced by $\underline{S_R^2}$ ($\because S_R^2 \xrightarrow{P} \sigma_R^2$)
 - probability of estimation error $\in [a, b]$, e.g., $P\left(\left|\frac{R - r_{xy}}{s_R}\right| > \delta\right) \approx 2[1 - \Phi(\delta)]$
 - approximate $100(1 - \alpha)\%$ confidence interval of \underline{r}_{xy} : $\underline{R} \pm z(\alpha/2) s_R$

10/2

Example 16 (estimate population ratio \underline{r}_{xy})

- Suppose that 100 people who recently bought houses are surveyed, and
 \underline{y} : mortgage payment \underline{x} : gross income
 are observed. The $\underline{r}_{xy} = \tau_y/\tau_x$ is the percentage of the total mortgage amount to the total gross income of all people who recently bought houses.
 - Suppose that the population size N is missing, but it is known that $100 \ll N$.
 - Suppose that $\bar{X} = 3100$, $\underline{s_x} = 1200$, $\bar{Y} = 868$, $s_y = 250$, $\hat{\rho}_{xy} = 0.85$. We have $\underline{R} = 868/3100 = 0.28$.
 - Neglecting the finite population correction, the estimated standard error of \underline{R} is $\underline{s_R} = \frac{1}{10} \times \frac{1}{3100} \sqrt{0.28^2(1200^2) + 250^2 - 2(0.28)(0.85)(250)(1200)} = 0.006$.
- Note that $\underline{s_R}$ is small because \underline{x} and \underline{y} are highly positively correlated, $\underline{r}_{xy} > 0$, and \bar{X} is large. ← check the graph in LNp.49

treated as S.R.S. with repl.

shows accuracy of R

without \approx with

- An approximate 95% confidence interval for \underline{r}_{xy} is
 $0.28 \pm 1.96 \times 0.006 = 0.28 \pm 0.012 = (0.268, 0.292)$. ← an interval estimate
- Again, neglecting the finite population correction, an estimated bias of \underline{R} using Thm 16 (LNp.48) is
 $\frac{1}{n} \times \frac{1}{\bar{X}^2} (Rs_x^2 - \hat{\rho}_{xy}s_x s_y) = \frac{1}{100} \times \frac{1}{3100^2} [(0.28)(250^2) - (0.85)(250)(1200)] = -0.00025$,
 which is negligible relative to $\underline{s_R}$ ($=0.006$). Note that the large $\hat{\rho}_{xy}$ ($=0.85$) and the large value of \bar{X} ($=3100$) cause the bias to be small.

• Ratios used for estimating population means (and totals)

- Suppose $\underline{\mu_x}$ is known, e.g., the example of 393 hospitals in Ex.2 (LNp.4),
 \underline{y} : number of discharges, ← main interest $\underline{\mu_y}$ ↑ population (N is known)
 \underline{x} : number of beds.

Why multiply \bar{Y} by this? Check the explanation in LNp.54

Suppose the average (or total) number of beds $\underline{\mu_x}$ (or τ_x) in the 393 hospitals is known (before a sample has been taken). Both $\underline{X_k}$ and $\underline{Y_k}$ are r.v.'s

- Q:** how to take advantage of this information in the estimation of $\underline{\mu_y}$?
- Select a random sample, and collect the data: $(\underline{X_k}, \underline{Y_k}), k = 1, \dots, n$. For the parameter $\underline{\mu_y} = \underline{\mu_x} \underline{r_{xy}}$, an intuitive ratio estimator of $\underline{\mu_y}$ is

$$\bar{Y}_R = \underline{\mu_x} \underline{R} = \bar{Y} \left(\frac{\underline{\mu_x}}{\bar{X}} \right) \quad (\leftarrow \text{cf. } \bar{Y}; \text{ Q: which estimator of } \underline{\mu_y} \text{ is better?}).$$

only use \underline{Y} data to estimate $\underline{\mu_y}$

use $(\underline{x}, \underline{y})$ data to estimate $\underline{\mu_y}$

- In the following discussion of this topic, we only consider the case of s.r.s. without replacement. The case of s.r.s. with replacement follows analogously. \rightarrow remove finite population correction

Example 17 (Comparison of sample mean and ratio estimator, cont. Ex.2 in LNp.4)

- Consider the example of hospital discharges. For the population of 393 (N) hospitals and $1 \leq i \leq N$, let **Observe Data** $(X_k, Y_k), k=1, \dots, n$ \rightarrow random, too.

positively highly correlated

- x_i : number of beds in the i th hospital (known before sampling)
- y_i : number of discharges in the i th hospital

- In this population,

\times : unknown in sampling survey parameters

$$\mu_x = 274.8 \text{ (known)},$$

$$\sigma_x = 213.2 \text{ (known)},$$

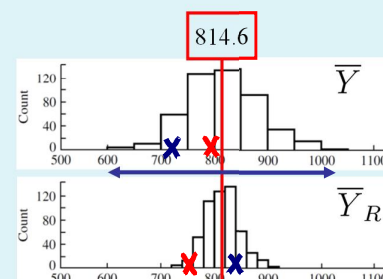
$$\mu_y = 814.6, \sigma_y = 589.7,$$

$$r_{xy} = 2.96, \rho_{xy} = 0.91.$$

For (X_k, Y_k) 's, same statistical modeling as given in LNp.41

Note. not

$$|\bar{Y} - \mu_y| > |\bar{Y}_R - \mu_y|$$



- To compare the performance of \bar{Y} and \bar{Y}_R , it was simulated (check Ex.3, LNp.15) 500 samples of size 64 (n) from the population of hospitals.

sampling distributions

$$MSE(\bar{Y}) > MSE(\bar{Y}_R)$$

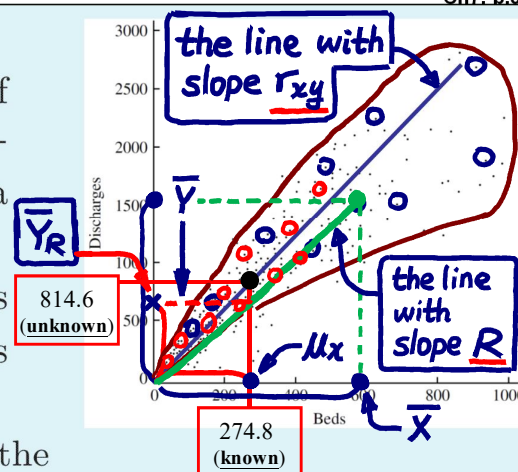
- The histograms of this result are shown in Figure 7.6 of textbook.

- The histograms show that the ratio estimator \bar{Y}_R of μ_y is less variable than the sample mean \bar{Y} .

- The comparison shows the ratio estimator \bar{Y}_R is effective at reducing variability $\Rightarrow \bar{Y}_R$ is a more accurate estimator than \bar{Y} . \leftarrow always true?

- Q:** Why is \bar{Y}_R better than \bar{Y} in this case?

- An explanation. Check the scatterplot of (x_i, y_i) for the 393 hospitals in the population (Figure 7.5 of textbook) and consider a random sample $(X_k, Y_k), k = 1, \dots, n$.



What if $\bar{X} < \mu_x$?

錨點

\downarrow

μ_x

μ_y

μ_x

μ_y

μ_x

μ_y

μ_x

μ_y

μ_x

μ_y

μ_x

μ_y

- the population correlation $\rho_{xy} = 0.91$ is high \Rightarrow a hospital with a large x_i tends to have a large y_i

- if $\bar{X} > \mu_x$, the sample over-estimates the number of beds μ_x , and probably the number of discharges as well, i.e., probably $\bar{Y} > \mu_y$.

- for this sample, multiplying \bar{Y} by $\frac{\mu_x}{\bar{X}}$ decreases \bar{Y} to \bar{Y}_R , which might be closer to μ_y than \bar{Y} .

$$R = \bar{Y}/\bar{X} = \bar{Y}_R/\mu_x$$

$$\frac{\mu_x}{\bar{X}} < 1$$

$$\bar{Y}(\mu_x/\bar{X})$$

Theorem 19 (approximate mean, bias, and variance of the ratio estimator)

Since $\bar{Y}_R = \mu_x R$, we have $E(\bar{Y}_R) = \mu_x E(R)$ and $Var(\bar{Y}_R) = \mu_x^2 Var(R)$. Under s.r.s. without replacement, Thm16 (LNp.48)

a known constant

- the approximate bias of the ratio estimator \bar{Y}_R of μ_y is

$$Bias(\bar{Y}_R)^2 \sim O(n^{-2})$$

$$E(\bar{Y}_R) - \mu_y = \mu_x [E(R) - r_{xy}] \approx \frac{1}{n} \left(1 - \frac{n-1}{N-1} \right) \times \frac{1}{\mu_x} (r_{xy} \sigma_x^2 - \rho_{xy} \sigma_x \sigma_y),$$

$$Bias(R) \uparrow \mu_y/\mu_x$$

not unbiased

Thm17 (LNp.49)

- The approximate variance of the ratio estimator \bar{Y}_R of μ_y is

$$\text{Var}(\bar{Y}_R) \sim O(n^{-1}) \quad \sigma_{\bar{Y}_R}^2 = \text{Var}(\bar{Y}_R) = \mu_x^2 \text{Var}(R) \approx \frac{1}{n} \left(1 - \frac{n-1}{N-1} \right) \times (r_{xy}^2 \sigma_x^2 + \sigma_y^2 - 2r_{xy} \rho_{xy} \sigma_x \sigma_y).$$

$\sigma_{xy} = \rho_{xy} \sigma_x \sigma_y$ ~~$\times \frac{1}{\mu_x^2}$~~

Proof: The results follows directly from the formulas for the approximate mean and variance of \bar{R} given in Thm. 16 (LNp.48) and Thm. 17 (LNp.49).

Note 13 (A note about the approximate variance of the ratio estimator)

Q: When will the ratio estimator \bar{Y}_R be better than the ordinary estimator \bar{Y} , i.e., $\text{Var}(\bar{Y}_R) < \text{Var}(\bar{Y})$? $\text{MSE} = \text{Var} + \text{Bias}^2$
 $O(n^{-1}) \sim \text{Var}$ $\sim O(n^{-2})$

- The ordinary estimator \bar{Y} has variance $\sigma_{\bar{Y}}^2 = \text{Var}(\bar{Y}) = \frac{\sigma_y^2}{n} \left(1 - \frac{n-1}{N-1} \right)$ (Thm. 3, LNp.18) and

$$\sigma_{\bar{Y}_R}^2 - \sigma_{\bar{Y}}^2 = \text{Var}(\bar{Y}_R) - \text{Var}(\bar{Y}) \approx \frac{1}{n} \left(1 - \frac{n-1}{N-1} \right) \times (r_{xy}^2 \sigma_x^2 - 2r_{xy} \rho_{xy} \sigma_x \sigma_y).$$

- The ratio estimator \bar{Y}_R has a smaller variance than \bar{Y} if

$$r_{xy}^2 \sigma_x^2 - 2r_{xy} \rho_{xy} \sigma_x \sigma_y < 0 \quad \Leftrightarrow \quad r_{xy}^2 \sigma_x < 2r_{xy} \rho_{xy} \sigma_y$$

$r_{xy} < 0 \leftarrow \rho_{xy} \rightarrow r_{xy} > 0$
-1 0 1

$$\text{i.e.} \quad \rho_{xy} \begin{cases} > \frac{1}{2} \left(\frac{\mu_y}{\mu_x} \right) \left(\frac{\sigma_x}{\sigma_y} \right) = \frac{1}{2} \left(\frac{CV_x}{CV_y} \right) > 0, & \text{provided that } r_{xy} > 0, \\ < \frac{1}{2} \left(\frac{\mu_y}{\mu_x} \right) \left(\frac{\sigma_x}{\sigma_y} \right) = \frac{1}{2} \left(\frac{CV_x}{CV_y} \right) < 0, & \text{provided that } r_{xy} < 0. \end{cases}$$

parameter
(free of unit)

where $CV_x = \sigma_x / \mu_x$ and $CV_y = \sigma_y / \mu_y$ are the coefficients of variation.