

Theorem 6 (variance of population total estimator)

- Under simple random sampling with replacement, $Var(\underline{T}) = N^2 \left(\frac{\sigma^2}{n} \right) \leftarrow Var(\bar{X})$
- Under simple random sampling without replacement,

cf. the precision of \bar{X}
 • Note 5 in LNp.17
 • Note 6 in LNp.19

$$Var(\underline{T}) = N^2 \left(\frac{\sigma^2}{n} \right) \left(1 - \frac{n-1}{N-1} \right) \leftarrow Var(\bar{X})$$

$$Var(\underline{T}) = Var(N\bar{X}) = N^2 Var(\bar{X})$$

Note. The precision of the estimator \underline{T} does depend on population size N .

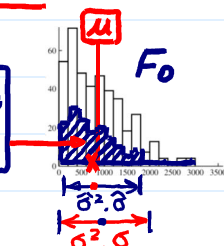
• Estimation of population variance

$$\sigma^2 = \sum_{j=1}^m \frac{n_j}{N} (C_j - \mu)^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

Recall. When F_0 is unknown, the σ in the standard error of \bar{X} is a parameter, i.e., it is unknown.

Q: how to estimate σ or σ^2 ?

histogram of data

**Definition 10 (sample variance)**

The sample variance of X_1, X_2, \dots, X_n is defined as $\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})^2$.

• a function of data • a r.v. • an estimator

Theorem 7 (expectation of sample variance, s.r.s. with replacement)

Under s.r.s. with replacement, we have $E(\hat{\sigma}^2) = \sigma^2 \left(\frac{n-1}{n} \right) \leftarrow \text{not unbiased}$

Proof. From the identity

$$\sum_{k=1}^n (X_k - \mu)^2 = \sum_{k=1}^n (X_k - \bar{X})^2 + n(\bar{X} - \mu)^2 \quad (\Delta)$$

by taking expectation on the both sides of (Δ) , we have

$$\sum_{k=1}^n E[(X_k - \mu)^2] = E \left[\sum_{k=1}^n (X_k - \bar{X})^2 \right] + n E[(\bar{X} - \mu)^2] \quad (\nabla)$$

which leads to

$$Var(X_k) \xrightarrow{=} \sigma^2 = E(\hat{\sigma}^2) + \sigma^2/n$$

Thus, we have $E(\hat{\sigma}^2) = ((n-1)\sigma^2)/n$.

Theorem 8 (expectation of sample variance, s.r.s. without replacement) not unbiased

Under s.r.s. without replacement, we have $E(\hat{\sigma}^2) = \sigma^2 \left(\frac{n-1}{n} \right) \left(\frac{N}{N-1} \right) \leftarrow$

Proof: The identities (Δ) and (∇) in the above proof still hold, and (∇) leads to

$$Var(X_k) \xrightarrow{=} \sigma^2 = E(\hat{\sigma}^2) + \sigma^2 \left(1 - \frac{1}{n} + \frac{n-1}{n} \times \frac{1}{N-1} \right)$$

After some algebra, this gives the desired result.

Note 7 (Some notes about the expectation of sample variance)

- No matter under s.r.s. with replacement or without replacement, the sample variance $\hat{\sigma}^2$ is a biased estimator of σ^2 .
- Since $\frac{n-1}{n} \leq 1$ and $\left(\frac{n-1}{n}\right) \left(\frac{N}{N-1}\right) < 1$ (note. $n < N$), we have

$$E(\hat{\sigma}^2) < \sigma^2.$$

That is, $\hat{\sigma}^2$ tends to underestimate σ^2 .

Theorem 9 (unbiased estimators of σ^2 and the variance of sample mean)

Under s.r.s. with replacement,

an unbiased estimator of σ^2 is

$$\frac{1}{n} \text{ (in } \hat{\sigma}^2) \approx \frac{1}{n-1} \text{ (in } S^2) \text{ when } n \text{ is large}$$

$$E(\hat{\theta}) = \theta \quad \text{a known constant}$$

$$E\left(\frac{\hat{\theta}}{a}\right) = \frac{\theta}{a}$$

$$s^2 = \left(\frac{n}{n-1}\right) \hat{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2$$

an unbiased estimator of $\text{Var}(\bar{X}) = \sigma^2/n$ is $s_{\bar{X}}^2 = s^2/n$.

although use same notation $S_{\bar{X}}^2$ different formula

Under s.r.s. without replacement,

an unbiased estimator of σ^2 is $\left(\frac{N-1}{N}\right) \left(\frac{n}{n-1}\right) \hat{\sigma}^2 = \left(\frac{N-1}{N}\right) s^2$

Thm 3 (LNp 18)

an unbiased estimator of $\text{Var}(\bar{X}) = (\sigma^2/n) \left(1 - \frac{n-1}{N-1}\right)$ is

sampling fraction

$$s_{\bar{X}}^2 = \frac{1}{n} \left(\frac{N-1}{N} s^2\right) \left(1 - \frac{n-1}{N-1}\right) = \frac{s^2}{n} \left(\frac{N-n}{N}\right) = \frac{s^2}{n} \left(1 - \frac{n}{N}\right)$$

Theorem 10 (unbiased est'ors of σ^2 and variance of sample mean, dichotomous x_i 's)

In the dichotomous cases, $\bar{X} = \hat{p}$ and $\sigma^2 = p(1-p)$.

Because

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})^2 = \frac{1}{n} \left(\sum_{k=1}^n X_k^2 \right) - \bar{X}^2 = \hat{p} - \hat{p}^2 = \hat{p}(1 - \hat{p})$$

we have,

unbiased est'or of population variance under with repl.

$$s^2 = \left(\frac{n}{n-1}\right) \hat{\sigma}^2 = \frac{n}{n-1} \hat{p}(1 - \hat{p})$$

functions of sample mean \hat{p} . When \hat{p} is known, $\hat{\sigma}^2, S^2$ are known.

Under s.r.s. with replacement, an unbiased estimator of $\text{Var}(\hat{p}) = \frac{p(1-p)}{n}$ is $s_{\hat{p}}^2 = s^2/n = [\hat{p}(1 - \hat{p})]/n - 1$.

Under s.r.s. without replacement, an unbiased estimator of $\text{Var}(\hat{p}) = \frac{p(1-p)}{n} \left(1 - \frac{n-1}{N-1}\right)$ is $s_{\hat{p}}^2 = \frac{s^2}{n} \left(1 - \frac{n}{N}\right) = \frac{\hat{p}(1 - \hat{p})}{n-1} \left(1 - \frac{n}{N}\right)$.

Theorem 11 (unbiased estimator of the variance of population total estimator)

An unbiased estimator of $\text{Var}(T) = N^2 \text{Var}(\bar{X})$ is $s_T^2 = N^2 s_{\bar{X}}^2$.

$$T = N\bar{X}$$

Thm 6 (LNp. 20)

unbiased estimator

$\frac{S^2}{n} (1 - \frac{n}{N})$ without
 $\frac{S^2}{n}$ with

The quantities $s_{\bar{X}}$ ($= \sqrt{s_{\bar{X}}^2}$), s_T ($= \sqrt{s_T^2}$), and $s_{\hat{p}}$ ($= \sqrt{s_{\hat{p}}^2}$) are called estimated standard errors. offer information about the magnitude of estimation error $\hat{\theta} - \theta$

Example 6 (estimate population mean, cont. Ex.2 in LNp.4)

An s.r.s. without replacement of size $n=50$ of the $N=393$ hospitals was taken.

discharge data:
 X_1, \dots, X_{50}

- From this sample, $\bar{X} = 938.5$ (recall, $\mu = 814.6$), $s = \sqrt{s^2} = 614.53$ (recall, $\sigma = 590$), and an estimate of $Var(\bar{X})$ is

population st.d.

$$s_{\bar{X}}^2 = \frac{s^2}{n} \left(1 - \frac{n}{N}\right) = \frac{614.53^2}{50} \left(1 - \frac{50}{393}\right) = 6592.$$

(unbiased) sample variance

sampling fraction 12.72%

- The estimated standard error of \bar{X} is $s_{\bar{X}} = \sqrt{6592} = 81.19$, $\rightarrow 2 \times s_{\bar{X}} = 162.38$

*: unknown in sampling survey

(cf. the (true) standard error of \bar{X} is $\sigma_{\bar{X}} = \frac{590}{\sqrt{50}} \sqrt{1 - \frac{49}{392}} = 78$) $\rightarrow E(\bar{X} - \mu)^2$ error

which gives a rough idea of how accurate the value of \bar{X} (938.5) is. In this case, the magnitude of the error is of the order 80, as opposed to 8 or 800.

- The error of \bar{X} is $938.5 - 814.9 = 123.9$, which is about $1.5 \times s_{\bar{X}}$. (cf.)

Example 7 (estimate population total, cont. Ex.2 in LNp.4)

- For the same sample in Ex.6, the estimate of the total number of discharges τ in the population of hospitals is $T = N \bar{X} = 393 \times 938.5 = 368,831$ (cf. the true value of τ is 320,139) \rightarrow population total

- The estimated standard error of T is $s_T = N s_{\bar{X}} = 393 \times 81.19 = 31,908$ (cf. the (true) standard error of T is $\sigma_T = N \sigma_{\bar{X}} = 393 \times 78 = 30,654$).

Example 8 (estimate population proportion, dichotomous x_i 's, cont. Ex.5 in LNp.21)

population proportion

- * $p = 0.654$: (true) proportion of hospitals in the population that had fewer than 1000 discharges ($\Rightarrow \sigma^2 = p(1-p) = 0.2263$). $\rightarrow \sum_{k=1}^{50} I_{[0,1000)}(X_k)$

population variance

- For the same sample in Ex.6 (LNp.26), 26 of 50 hospitals has fewer than 1000 discharges. The estimate of p is $\hat{p} = 26/50 = 0.52$, and an estimate of $Var(\hat{p})$ is

sample proportion

*: unknown in sampling survey

(unbiased) sample variance

$$s_{\hat{p}}^2 = \frac{\hat{p}(1-\hat{p})}{n-1} \left(1 - \frac{n}{N}\right) = \frac{(.52)(.48)}{49} \left(1 - \frac{50}{393}\right) = 0.0045$$

- The estimated standard error of \hat{p} is $s_{\hat{p}} = \sqrt{0.0045} = 0.067$, $\rightarrow 2 \times s_{\hat{p}} = 0.134$

$E(\hat{p} - p)^2$ error

(cf. (true) standard error of \hat{p} is $\sigma_{\hat{p}} = \sqrt{\frac{(.654)(.346)}{50}} \sqrt{1 - \frac{49}{392}} = 0.064$)

which tells us that the error of \hat{p} is in the 2nd or 1st decimal place — we are probably not so fortunate as to have an error in the 3rd decimal place.

- The true error of \hat{p} is $0.52 - 0.654 = -0.134$, which is about $-2 \times s_{\hat{p}}$. (cf.)

- Note. Examples 6-8 show how, in s.r.s., we can not only form estimates of unknown population parameters (e.g., use \bar{X} , T , \hat{p} to estimate μ , τ , p , respectively), but also gauge the likely size of the errors of the estimates, by estimating their standard errors (e.g., $s_{\bar{X}}$, s_T , $s_{\hat{p}}$) using the data in the sample.

Note 8 (A summary of parameter estimation in s.r.s.)

- A summary table: **parameter** **statistic (random variable)**

population parameter	estimator ^(†)	variance of estimator ^{(†)(*)}	estimated variance ^{(†)(*)}
μ	$\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$	(a) $\frac{\sigma_X^2}{n} = \frac{\sigma^2}{n}$ (b) $\frac{\sigma_X^2}{n} = \frac{\sigma^2}{n} \left(1 - \frac{n-1}{N-1}\right)$	(a) $\frac{s_X^2}{n} = \frac{s^2}{n}$ (b) $\frac{s_X^2}{n} = \frac{s^2}{n} \left(1 - \frac{n}{N}\right)$
p	$\hat{p} = \text{sample proportion}$	(a) $\frac{\sigma_p^2}{n} = \frac{p(1-p)}{n}$ (b) $\frac{\sigma_p^2}{n} = \frac{p(1-p)}{n} \left(1 - \frac{n-1}{N-1}\right)$	(a) $\frac{s_p^2}{n} = \frac{\hat{p}(1-\hat{p})}{n-1}$ (b) $\frac{s_p^2}{n} = \frac{\hat{p}(1-\hat{p})}{n-1} \left(1 - \frac{n}{N}\right)$
τ	$T = N \bar{X}$	(a) $\sigma_T^2 = N^2 \sigma_X^2$ (b) $\sigma_T^2 = N^2 \sigma_X^2$	(a) $s_T^2 = N^2 s_X^2$ (b) $s_T^2 = N^2 s_X^2$
σ^2	(a) $s^2 = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X})^2$ (b) $\left(1 - \frac{1}{N}\right) s^2$	$\frac{n\hat{p}(1-\hat{p})}{n-1}$ dichotomous case	$\frac{n\hat{p}(1-\hat{p})}{n-1}$ dichotomous case

When $n \ll N$, almost identical

finite population correction

related to sampling fraction n/N

- (†): (a) and (b) obtained under with and without replacement, respectively.
- (*): the square root of entries in the 3rd column are standard errors, the square root of entries in the 4th column are estimated standard errors.

