## • Normal approximation to the sampling distribution of sample mean

*Recall. the shape of $F_{\overline{X}}$ in LNp.14~15*

**We have known**
- $E(\overline{x}) = \mu$
- $Var(\overline{x}) \to 0$, $n \to \infty$

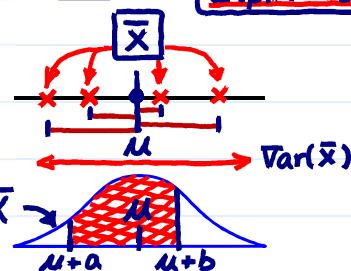*how about the shape of $F_{\overline{x}}$?*

*Recall, dichotomous case (Thm4, LNp 20)*

- **Q**: without knowledge of the population distribution $F_0$, ←unknown how to further characterize the sampling distribution $F_{\overline{X}}$ ←shape = ? of $\overline{X}$ in addition to its mean and variance?

- Advantages if we (almost) know the shape of $F_{\overline{X}}$?
  - accurately evaluate $P(\text{error} \in (a,b)) \approx \underline{?}$
  - (**Note.** error $= \overline{X} - \mu$)
  - construct confidence interval for $\mu$

*pdf/pmf of $\overline{X}$*

$(\bigstar)$ $\dfrac{\overline{X}_n - \mu}{S_{\overline{x}}} \overset{D}{\approx} N(0,1)$

$\boxed{\dfrac{\overline{X}_n - \mu}{\sigma_{\overline{x}}} \overset{D}{\approx} N(0,1)}$

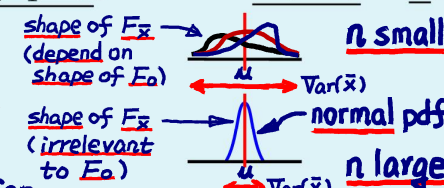| **Theorem 12** (central limit theorem, CLT, for i.i.d. case) ← S.R.S with repl. |
| --- |

Suppose that $X_1, X_2, \ldots, X_n$ are i.i.d. r.v.'s and have common mean $\mu$ and variance $0 < \sigma^2 < \infty$. For the sample mean $\overline{X}_n = \frac{1}{n}\sum_{k=1}^{n} X_k$, we have $E(\overline{X}_n) = \mu$, $\sigma_{\overline{X}_n}^2 = Var(\overline{X}_n) = \sigma^2/n$, and for any fixed value $z$,

*cdf of $\frac{\overline{X}_n - \mu}{\sigma/\sqrt{n}}$*

$$P\left(\frac{\overline{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{\overline{X}_n - \mu}{\sigma_{\overline{X}_n}} \le z\right) \longrightarrow \Phi(z)$$

*standardization*

*shape of $F_{\overline{x}}$ (depend on shape of $F_0$)* — n small
*shape of $F_{\overline{x}}$ (irrelevant to $F_0$)* — normal pdf, n large

*Then. mean=0, var=1*

as $n \to \infty$, where $\Phi$ is the cumulative distribution function (cdf) of the standard normal distribution $N(0,1)$. That is, $\overline{X}_n \overset{D}{\approx} N(\mu, \sigma^2/n)$. *Shape of $F_{\overline{x}}$*

(**cf.**) Law of large number (LLN) guarantees that $\overline{X}_n \overset{P}{\longrightarrow} \mu$ and $s^2 \overset{P}{\longrightarrow} \sigma^2$ *TBp.179 ~180* as $n \to \infty$, i.e., $\overline{X}_n$ and $s^2$ are **consistent** estimators of $\mu$ and $\sigma^2$, respectively.

| **Theorem 13** (central limit theorem, CLT, for s.r.s. without replacement) |
| --- |

In s.r.s. without replacement, (1) $X_1, X_2, \ldots, X_n$ are not independent, and (2) there is no reason to have $n \to \infty$ while $N$ remains fixed. But other CLTs are still appropriate, e.g.,

$0 \ll n \ll N$

If $n$ is large, but still small relative to $N$, CLT← without $\approx$ with

then $\overline{X}_n$ is approximately normally distributed with mean $\mu$ and variance $\sigma_{\overline{X}_n}$ $= (\sigma/\sqrt{n}) \cdot \sqrt{1 - \frac{n-1}{N-1}}$, not $\sigma/\sqrt{n}$ (check graphs in Ex.3, LNp.15).

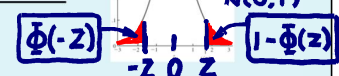| **Application 1** (CLT application on estimation error of population mean) |
| --- |

A use of CLT for estimation error $\overline{X}_n - \mu$ is

*$\overset{D}{\approx} N(0,1)$*

$\boxed{\sigma_{\overline{x}} \text{ as a standard}}$

$$P\left(|\overline{X}_n - \mu| < \delta\right) = P\left(-\delta \le \overline{X}_n - \mu \le \delta\right) = P\left(-\frac{\delta}{\sigma_{\overline{X}_n}} \le \frac{\overline{X}_n - \mu}{\sigma_{\overline{X}_n}} \le \frac{\delta}{\sigma_{\overline{X}_n}}\right)$$

*|error|*    $1 - \Phi(\delta/\sigma_{\overline{x}})$    *usually unknown*

$$\approx \Phi\left(\frac{\delta}{\sigma_{\overline{X}_n}}\right) - \Phi\left(-\frac{\delta}{\sigma_{\overline{X}_n}}\right) = 2\,\Phi\left(\frac{\delta}{\sigma_{\overline{X}_n}}\right) - 1.$$

*pdf of $N(0,1)$*

| $\pm\delta$ | $\pm 2\sigma_{\overline{x}}$ | $\pm 1.96\sigma_{\overline{x}}$ | $\pm 1.64\sigma_{\overline{x}}$ | $\pm 1\sigma_{\overline{x}}$ |
| --- | --- | --- | --- | --- |
| $\Phi$ | 0.954 | 0.950 | 0.900 | 0.685 |

- **Note.** For the cdf $\Phi$ of $N(0,1)$, $\Phi(-z) = 1 - \Phi(z)$.

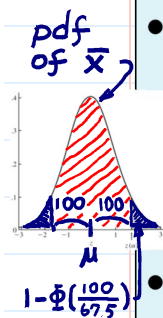$\boxed{\Phi(-z)}$    $\boxed{1-\Phi(z)}$    $-z \; 0 \; z$

| **Example 9** (probability of estimation error more than δ, cont. Ex.2 in LNp.4) |
| --- |

- Consider the population of 393 hospitals and s.r.s. without replacement.

- For $n = 64$, $\overset{\times}{\sigma_{\overline{X}}} = \dfrac{\overset{\times}{\sigma}}{\sqrt{n}}\sqrt{1 - \dfrac{n-1}{N-1}} = \dfrac{\overset{\times}{589.7}}{8}\sqrt{1 - \dfrac{63}{392}} = \overset{\times}{67.5}$.

*×: unknown in sampling survey*

pdf of $\overline{X}$

$1-\Phi\left(\frac{100}{67.5}\right)$

- Apply <u>CLT</u> to <u>approximate</u> the <u>probability</u> that the <u>sample mean</u> $\overline{X}$ differs from $\mu$ by <u>more than</u> $\delta = 100$:

  $N(0,1) \overset{D}{\approx}$

  $$P(|\overline{X} - \mu| > \underline{100}) = \underline{2} \times P(\overline{X} - \mu > \underline{100}) = 2 \times P\left(\frac{\overline{X} - \mu}{\times \sigma_{\overline{X}}} > \frac{100}{\times \sigma_{\overline{X}}}\right)$$

  $N(0, \sigma_{\overline{x}}^2)$    $\approx$    $2\left[1 - \Phi\left(100/67.5\right)\right] = 2 \times \underline{0.069} = \underline{0.14}.$   cf.

  ×: unknown in sampling survey

- Among <u>500 samples</u> of <u>size 64</u> (Ex.3, LNp.15), <u>82 samples</u> (or <u>16.4%</u>) <u>differed</u> from $\mu$ more than <u>100</u>.   $n\hat{p} \sim$ binomial (**with**) or hypergeometric (**without**)

  Thm 4 LNp.20

**Example 10** (estimation error <u>more than</u> $\delta$, <u>dichotomous</u> $x_i$'s, cont. <u>Ex.8 in LNp.27</u>)

pdf of $\hat{p}$

$1-\Phi(2.094)$

- <u>sample size</u> $n = 50$, true $\overset{\times}{p} = 0.654$, standard error of $\hat{p}$ is $\overset{\times}{\sigma_{\hat{p}}} = 0.064.$   estimate

- From the <u>sample</u> in <u>Ex.8</u>, estimate of $p$ is $\hat{p} = 0.52$ and $|\hat{p} - p| = 0.134$, the <u>probability</u> that the estimator will be <u>off</u> by an amount <u>this large or larger</u> is

  estimator $\rightarrow$ $P(|\hat{p} - p| > 0.134) = \underline{1 - P}(|\hat{p} - p| \leq 0.134)$   Recall. Normal approximation to binomial (TBp.187)

  $N(0, \sigma_{\hat{p}}^2)$   $= 1 - P\left(\frac{|\hat{p} - p|}{\times \sigma_{\hat{p}}} \leq \frac{0.134}{\times \sigma_{\hat{p}}}\right) \approx 2\left[1 - \Phi(2.094)\right] = \underline{0.036}.$

  $N(0,1) \overset{D}{\approx}$

  ×: unknown in sampling survey

- We see that the <u>sample</u> was rather "<u>unlucky</u>" — an error <u>this large or larger</u> would <u>occur</u> only about $\overset{\times}{3.6\%}$ of the time.

(★) in LNp.29

**Note.** In a <u>sampling survey</u>, $\sigma^2$ (or $\sigma_{\overline{X}_n}^2$) is <u>not</u> available because $F_0$ remains <u>unknown</u>. We can <u>substitue</u> $s^2$ for $\sigma^2$,      cdf of $(\overline{X} - \mu)/s_{\overline{x}}$
and a <u>similar CLT still holds</u>, i.e.,   $P\left(\frac{\overline{X}_n - \mu}{s_{\overline{X}_n}} < \underline{z}\right) \longrightarrow \Phi(z)$   as $\underline{n \to \infty}$.

$N(0,1) \overset{D}{\approx}$   [(a) with (b) without]

---

**Definition 11** (<u>interval estimator</u>, <u>coverage probability</u>, <u>interval estimate</u>, <u>confidence interval</u>, and <u>confidence level</u>)   ← Why construct confidence interval?
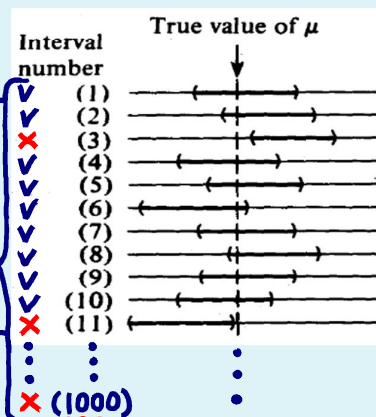
Data

- For a <u>random vector</u> $\mathbf{X} = (X_1, \ldots, X_n)$, an **interval estimator** of a parameter $\theta$ with **coverage probability** $1 - \alpha$ is a <u>random interval</u>

  $$(\hat{\theta}_L(\mathbf{X}), \hat{\theta}_U(\mathbf{X})),$$

  where    statistics (r.v.'s)

  Repeated construction of **95%** confidence intervals

$\theta$ is contained in the interval estimator

  1. $\hat{\theta}_L(\mathbf{X}), \hat{\theta}_U(\mathbf{X})$ are <u>functions of data only</u>,
  2. $\hat{\theta}_L(\mathbf{X}) < \hat{\theta}_U(\mathbf{X})$, and,   not-covered probability
  3. $\underline{P(\theta \in (\hat{\theta}_L(\mathbf{X}), \hat{\theta}_U(\mathbf{X})))} = \underline{1 - \alpha}.$

  about 950 interval estimates containing $\mu$



| Interval number | True value of $\mu$ |
|---|---|
| ✓ (1) | |
| ✓ (2) | |
| ✗ (3) | |
| ✓ (4) | |
| ✓ (5) | |
| ✓ (6) | |
| ✓ (7) | |
| ✓ (8) | |
| ✓ (9) | |
| ✓ (10) | |
| ✗ (11) | |
| ⋮ | |
| ✗ (1000) | |

- If $\mathbf{X} = \mathbf{x}$ is observed, the interval

observed data

  $$(\hat{\theta}_L(\mathbf{x}), \hat{\theta}_U(\mathbf{x}))$$

  either contains $\mu$ or not, i.e., probability of containing $\mu$ is 1 or 0

  is called an **interval estimate**.

- The term "$100 \times (1 - \alpha)\%$ **confidence interval**" (**C.I.**) is used to denote either an interval estimator with coverage probability $1 - \alpha$ or an interval estimate.
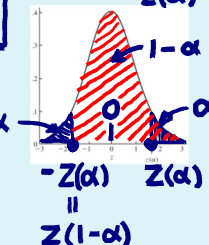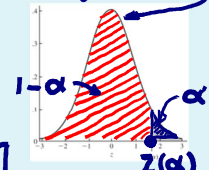
  significance level $\alpha$ in testing

- The $100(1 - \alpha)\%$ is also referred to as **confidence level**.   cf.

  90% ←   → 95%   → 99%

- **Note.** The $\alpha$ is usually assigned a <u>small</u> value, e.g. <u>0.1</u>, <u>0.05</u>, or <u>0.01</u>.

## Application 2 (CLT application on the construction of confidence interval for μ)

- For $0 \leq \alpha \leq 1$, let $z(\alpha)$ be the $(1-\alpha)$-quantile of $N(0,1)$, i.e., $z(\alpha)$ is the number such that the area under the pdf of $N(0,1)$ to the right of $z(\alpha)$ is $\alpha$ and $\Phi(z(\alpha)) = 1-\alpha$. **Notice** that $z(1-\alpha) = -z(\alpha)$.

*quantile 分位數*

*pdf of N(0,1)*

*1-α*   *α*   *z(α)*

- For $Z \sim N(0,1)$, $P\big(-z(\alpha/2) \leq Z \leq z(\alpha/2)\big) = \Phi(z(\alpha/2)) -$

$$\underset{1-\Phi(z(\alpha/2))}{} \quad \Phi(-z(\alpha/2)) = 2 \times \Phi(z(\alpha/2)) - 1 = 1-\alpha.$$

*1-α/2*

*We know this so that we can construct C.I.*

*1-α*   *α*   *α*

*-Z(α)*   *0*   *Z(α)*

*= Z(1-α)*

- Because $\overline{X}_n \overset{D}{\approx} N(\mu, \sigma^2_{\overline{X}_n})$ by CLT, we have

*usually unknown, assume known here*

$$* \frac{\overline{X}-\mu}{S_{\overline{X}}} \overset{D}{\approx} N(0,1)$$ *by (\*) in LNp. 29*

*(asymptotic) pivotal quantity of μ*   *if σ² known*

$$P\Big(-z(\alpha/2) \leq \frac{\overline{X}_n - \mu}{\sigma_{\overline{X}_n}} \leq z(\alpha/2)\Big) \overset{\approx}{} 1-\alpha$$

$$\overset{D}{\approx} N(0,1)$$

$$\Leftrightarrow P\Big(\overline{X}_n - z(\alpha/2)\sigma_{\overline{X}_n} \leq \mu \leq \overline{X}_n + z(\alpha/2)\sigma_{\overline{X}_n}\Big) \approx 1-\alpha$$

- The probability that $\mu$ lies in the random interval formed by data:

*check Def. 11, 1,2,3 in LNp.32* →

$$\overline{X}_n \pm z(\alpha/2)\sigma_{\overline{X}_n}$$

*replaced by $S_{\overline{X}}$*

*x: unknown in sampling survey*

is $\approx 1-\alpha$, i.e., it is a $100(1-\alpha)\%$ (asymptotic) confidence interval of $\mu$.

- **Recall**. A function $Q(\mathbf{X}, \theta)$ of the data $\mathbf{X}$ and a parameter, say $\theta$, of interest is called a **pivotal quantity** for $\theta$ if the distribution of $Q(\mathbf{X}, \theta)$ is irrelevant to all parameters.

*a r.v., but not a statistic →*

## Note 9 (Some notes about confidence interval)

- In a sample survey, $\sigma_{\overline{X}_n}$ is unknown. In the case, $s_{\overline{X}_n}$ (or $s^2$, respectively) can be substituted for $\sigma_{\overline{X}_n}$ (or $\sigma^2$, respectively) if the sample size $n$ is large enough, say $n \geq 25$ or $30$ by a rule of thumb.

*TBp.338*

- **Recall**: duality between confidence interval and hypothesis testing.
  - Suppose for every parameter value $\theta_0$, there is a level-$\alpha$ test for
  $$H_0 : \theta = \theta_0 \quad \text{vs.} \quad H_A : \theta \neq \theta_0.$$
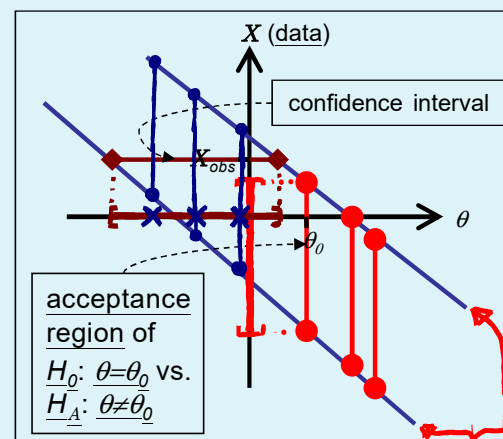  Denote the acceptance region of the test by $AR(\theta_0)$. Then, the set
  $$C(\mathbf{X}) = \{\theta \,|\, \mathbf{X} \in AR(\theta)\}$$
  is a $100(1-\alpha)\%$ C.I. for $\theta$.

*a set of parameter values*

*X (data)*

*confidence interval*

*acceptance region of $H_0$: $\theta = \theta_0$ vs. $H_A$: $\theta \neq \theta_0$*

*$X_{obs}$*   *$\theta_0$*   *$\theta$*

  - Suppose $C(\mathbf{X})$ is a $100(1-\alpha)\%$ C.I. for $\theta$. Then, an acceptance region for a level-$\alpha$ test of $H_0$: $\theta = \theta_0$ is
  $$AR(\theta_0) = \{\mathbf{X} \,|\, \theta_0 \in C(\mathbf{X})\}.$$

*can use C.I. to perform a test*

- In a sample survey, for the population mean $\mu$ and the hypotheses $H_0$: $\mu = \mu_0$ vs. $H_A$: $\mu \neq \mu_0$, a test at (aymptotic) significance level $\alpha$ rejects $H_0$ if

*distance difference*

*known value*

$$\Big|(\overline{X}_n - \mu_0)/\sigma_{\overline{X}_n}\Big| > z(\alpha/2)$$

*pivotal quantity →* $\Big|\frac{\overline{X}-\mu}{\sigma_{\overline{X}}}\Big| < z(\frac{\alpha}{2})$ *in LNp 33*

*cf.*

*scale marked on a ruler*

- Many <u>confidence intervals</u> have the <u>form</u>:

$$\underline{\text{estimate}} \pm \underline{[\text{critical value}]} \times \underline{[(\text{estimated})\ \text{standard error}]}$$

*(margin, top right: $\hat{\theta}$)*

$\Rightarrow$ C.I. <u>combines</u> <u>information</u> of <u>estimate</u> and (estimated) <u>standard error</u>

- The <u>width</u> of a <u>confidence interval</u> <u>often</u> <u>depends</u> <u>on</u>:

<u>Under same $\alpha$, Smaller width $\Leftrightarrow$ more accurate C.I.</u>

    –  $\underline{n}$ : <u>sample size</u>

        $\underline{n \uparrow}$,  <u>width $\downarrow$</u>

    –  $\underline{\sigma}$: <u>population standard deviation</u>

        $\underline{\sigma \uparrow}$,  <u>width $\uparrow$</u>

    –  $\underline{1-\alpha}$: <u>confidence level</u>

        $\underline{(1-\alpha) \uparrow}$,  <u>width $\uparrow$</u>

*(margin: e.g., use previously collected information of population)*

For example, consider the <u>C.I.</u>:

$$\overline{X}_n \pm \underline{z(\alpha/2)} \times \underline{\sigma_{\overline{X}_n}}$$
$$= \ \overline{X}_n \pm \underline{z(\alpha/2)} \times \frac{\sigma}{\sqrt{n}}$$

*(margin: $s \leftarrow$ / $\sigma$)*

- If $\alpha$ is <u>fixed</u> and $\underline{\sigma}$ is (approximately) <u>known</u>, $\underline{n}$ can be <u>chosen</u> so as to obtain <u>confidence intervals</u> close to some <u>desired length</u>. $\rightarrow$ *i.e., use estimated st.e. to determine $n$*

$\Rightarrow$ a <u>common way</u> to determine an <u>adequate</u> survey <u>sample size</u> $n$

---

**Example 11** (<u>repeated construction</u> of <u>confidence intervals</u>, cont. <u>Ex.2</u> in <u>LNp.4</u>)

- <u>20 samples</u> each of <u>size</u> $n=25$ were <u>drawn</u> from the <u>population</u> of <u>hospital discharges</u> $(N=393)$.

- From <u>each</u> of the <u>samples</u>, an (approximate) <u>95%</u> <u>confidence interval</u> for $\underline{\mu}$ was <u>computed</u> and <u>displayed</u> in <u>Figure 7.4</u> (textbook).

- On <u>average</u> <u>5%</u>, or <u>1</u> out of <u>20</u>, would <u>not</u> <u>include</u> $\underline{\mu}$.



*(margin near figure: $\mu$   814.6    all 20 C.I.'s contain $\mu$)*

---

**Example 12** (<u>construction</u> of <u>confidence intervals</u> for $\underline{\mu}$, $\underline{\tau}$, $\underline{p}$)

- A <u>particular area</u> contains <u>8000</u> (population size $\underline{N}$) condominium <u>units</u>.

- To <u>understand</u> the <u>numbers</u> of motor <u>vehecles</u> owned by the <u>units</u>, a s.r.s. <u>without</u> replacement of size $n=100$ was <u>drawn</u>. $\leftarrow x_i, i=1, \cdots, 8000$

- The <u>sample yields</u> that  $\rightarrow$ *Data: $X_1, \cdots, X_{100}$*     *parameter $\rightarrow$* / *estimate $\rightarrow \mu$*

    – the <u>average number</u> of <u>motor vehicles</u> per <u>unit</u> is $\overline{X}=1.6$,

    – with a <u>sample standard deviation</u> $s=0.8$. $\leftarrow s^2$ *estimates population variance*

    – So,

$$s_{\overline{X}} = \frac{s}{\sqrt{n}}\sqrt{1-\frac{n}{N}} = \frac{0.8}{\sqrt{100}}\sqrt{1-\frac{100}{8000}} = \underline{0.08}. \rightarrow \text{shows accuracy of } \overline{X}$$

*(margin: $s_{\overline{X}}^2$ estimates the variance of $\overline{X}$)*

- When $\alpha=0.05$, we have $z(\alpha/2)=z(0.025)=\underline{1.96}$. Therefore, a <u>95%</u> <u>confidence interval</u> for the <u>population average</u> $\mu$ is

*(box: C.I. combines 2 information)*

$$\overline{X} \pm 1.96 \times s_{\overline{X}} = (1.44, 1.76).$$

*Why?*    *an interval estimate: a collection of many possible $\mu$'s*

- For the <u>population total</u> $\underline{\tau} = \underline{N\mu}$ (i.e., <u>total number</u> of <u>motor vehicles</u> owned by the <u>8000 units</u>), $\llcorner$ *parameter*

    – an <u>estimate</u> of $\underline{\tau}$ is $\underline{T} = \underline{N \times \overline{X}} = 8000 \times 1.6 = \underline{12,800}$,

    – with an <u>estimated standard error</u> $s_T = \underline{N \times s_{\overline{X}}} = \underline{640}. \rightarrow$ *shows accuracy of $T$*

- So, a <u>95%</u> <u>confidence interval</u> for $\tau$ is   *8000×(1.44,1.76)*

$$\underline{T \pm 1.96 \times s_T} = (11,546, 14,054). \leftarrow \text{an interval estimate}$$

*Why?*

- In the sample, 12% of the 100 $(n)$ respondents said that they plan to sell their condos within the next year. → *dichotomous data*

- For the proportion $p$ of 8000 $(N)$ units whose owners were planning to sell the units in next year, ↳ *parameter*

  *shows accuracy of $\hat{p}$* ←

  – an estimate of $p$ is $\hat{p} = 0.12$,

  8000 ×
  (0.06, 0.18)

  – with an estimated standard error $s_{\hat{p}} = \sqrt{\dfrac{\hat{p}(1-\hat{p})}{n-1}}\sqrt{1-\dfrac{n}{N}} = 0.03.$

- So, a 95% confidence interval for $p$ is $\hat{p} \pm 1.96 \times s_{\hat{p}} = (0.06, 0.18).$ ← *interval estimate*

- ⊙ A 95% confidence interval for the total number $(=N\times p)$ of owners planning to sell is

$$(N\hat{p}) \pm 1.96 \times (N s_{\hat{p}}) = (451, 1469).$$ ← *if this too wide*

**Example 13** (sample size determination, cont. Ex.12 in LNp.36)

$n=100$

509
⌢
$N\hat{p}$

$n^* = ?\, (>n)$

200
⌢
$N\hat{p}$

- Suppose a 95% C.I. of $Np$ with a half-width of 200 is desired (cf., original half-width: $(1469 - 451)/2 = 509$).

  *why can we do this?* ⇒ *assume $n^* \ll N$*

- For a sample of size $n^*$, half-width of 95% C.I. of $Np$, neglecting the finite population correction (i.e., treated as s.r.s. with replacement), is

*Why use s.r.s. with?* ⇒ *easier to solve for $n^*$*

$$1.96 \times (N s_{\hat{p}}) \underset{\approx}{\approx} 1.96 \times N\sqrt{\dfrac{\hat{p}(1-\hat{p})}{n^*}} = \dfrac{5095}{\sqrt{n^*}} \cdot$$

*previously collected information*

$\dfrac{649}{100}$
$= 6.51\frac{\#}{\Xi}$
$\approx \left(\dfrac{509}{200}\right)^2$
∵ $S_{\hat{p}}$
∝ $1/\sqrt{n}$

*S.r.s. with*

- Setting $5095/\sqrt{n^*} = 200$ and solving for $n^*$, we have $n^* = (5095/200)^2 = 649$ (cf., original sample size $n$: 100).

  ↳ $\ll N \approx 8000$

❖ **Reading**: textbook, 7.3