

Here we want to draw conclusions about the pop from the sample

Unbiased Estimates: $E(f_{\bar{x}}) = E(f) = \mu$ (centered at same value)
 so \bar{x} is an unbiased estimator of μ . a.k.a. $\hat{\mu} = \bar{x}$

Efficient Estimates: All the sampling distrib of 2 statistics have the same Expected value (center) then the one with the smaller var. is more efficient.

Consistent Estimate: As sample size increases, the statistic converges in probability to pop value.

Confidence Interval for μ

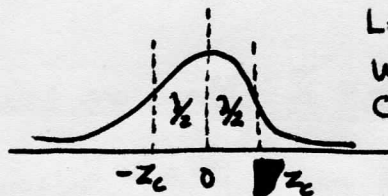
We know $f_{\bar{x}} \approx f_{N(\mu, \sigma_{\bar{x}}^2)}$ thus we do calculations with N distrib.

How do we derive a confidence interval for μ ? Say we want 95% CI

Let $\lambda = 0.95$.

We want the area of the center strip in $N(0,1)$ to be λ .

By symm, there are 2 identical areas of size $\lambda/2$.



[If $\lambda = 0.95$, we want area 0.4750 on each side, so we seek $z_c \ni Q(z_c) = 0.4750 \Rightarrow z_c = 1.96$
 ↑ value we look up in table.]

$$Z = \frac{\bar{x} - E(\bar{x})}{\sqrt{\text{var}(\bar{x})}} = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$$

so we'd have

$$\begin{aligned} -z_c &< \frac{\bar{x} - \mu}{\sigma_{\bar{x}}} < z_c \\ -z_c \sigma_{\bar{x}} &< \bar{x} - \mu < z_c \sigma_{\bar{x}} \\ -z_c \sigma_{\bar{x}} - \bar{x} &< -\mu < z_c \sigma_{\bar{x}} - \bar{x} \end{aligned}$$

$$\Rightarrow \bar{x} - z_c \sigma_{\bar{x}} < \mu < \bar{x} + z_c \sigma_{\bar{x}}$$

$$\Rightarrow \mu \in (\bar{x} - z_c \sigma_{\bar{x}}, \bar{x} + z_c \sigma_{\bar{x}})$$

Here we can substitute $\sigma_{\bar{x}} = \frac{1}{\sqrt{N}} \sigma$ or $\sqrt{\frac{N-1}{N}} \frac{1}{\sqrt{N}} \sigma$

Common values for z_c

95%	1.96
99%	2.58
99.97%	3.0

This is interpreted to mean $\mu \in (\bar{x} - 1.96\sigma, \bar{x} + 1.96\sigma)$ with Prob = 0.95.

Objections: Neyman said it was only true that if you computed many samples \bar{x} and formed these intervals that 95% of the time μ would be in an interval. He is saying that once a specific $\{x_1, \dots, x_n\}$ are chosen, \bar{x} is determined and there is no randomness - either $\mu \in CI$ or $\mu \notin CI$ with certainty (Prob = 1).
 (This is like saying before I flip a coin, $P(\text{heads}) = 1/2$ but after I flip it, (even if I don't look) $P(\text{heads}) = 1$ or 0 .)

E. Lehmann
 'Fisher, Neyman and Creation of Classical Stats' p. 89

I don't see this as a problem. We compute \bar{x} and form CI and there is a 95% chance that this is a good \bar{x} and thus $\mu \in CI$.
 The only problem I see is that this holds for any \bar{x} I compute. \bar{x}_1 and \bar{x}_2 may have very different values, but there is no way to distinguish which is better (unless we repeatedly do it and try to compute $f_{\bar{x}}$ and see its center!)

⑧ Given $\sigma = 0.05$, how large must sample size N be to make the 95% CI bounds ≤ 0.01 ? How about for 99% CI?

$$\bar{X} \pm z_c \frac{\sigma}{\sqrt{N}} \text{ so we want } z_c \frac{\sigma}{\sqrt{N}} \leq 0.01 \Rightarrow N \geq \left(\frac{z_c \sigma}{0.01} \right)^2$$

For 95% $z_c = 1.96$ $N \geq \left(\frac{1.96 (0.05)}{0.01} \right)^2 = 96.04 \Rightarrow N > 97$

99% $z_c = 2.58$ $N \geq \left(\frac{2.58 (0.05)}{0.01} \right)^2 = 166.41 \Rightarrow N > 167$

⑩ In a sample of $N=100$, 55 people said they would vote for the candidate. What is a 95% CI for the proportion of the pop who will vote for the candidate?

From ch 8, we know proportions are really $\bar{X} = \frac{1}{N} \sum X_i$; $X_i = \begin{cases} 1 \\ 0 \end{cases}$

Estimate of population prob from sample: $\hat{p} = 0.55$ we have $\hat{p} = \bar{x}$
 $\mu_{\hat{p}} = \mu_{\bar{x}} = \mu$
 $\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{N}}$

$$-z_c \leq \frac{\hat{p} - \mu_p}{\sigma_p} \leq z_c$$

$\Rightarrow \hat{p} - z_c \sigma_p < \mu_p < \hat{p} + z_c \sigma_p$ we have to use \hat{p} to estimate σ_p

$\mu_p \in \left(0.55 - \frac{1.96 (0.0497)}{0.1}, 0.55 + \frac{1.96 (0.0497)}{0.1} \right)$ $\sigma_p = \sqrt{\frac{(0.55)(0.45)}{100}} = 0.0497$

For 99% CI,
 $\mu_p \in \left(0.55 - \frac{2.58 (0.0497)}{0.1}, 0.55 + \frac{2.58 (0.0497)}{0.1} \right)$

⑪ For prob ⑩, how large a sample size N should we take so that we have { 95% CI that the candidate will get elected? }
 { 99% }

To win, the candidate must receive $> 50\%$ of vote in pop.

$$\hat{p} - z_c \sigma_p = \hat{p} - z_c \frac{\sqrt{\hat{p}(1-\hat{p})}}{\sqrt{N}} \stackrel{!}{>} 0.50 \quad \sqrt{\hat{p}(1-\hat{p})} = \sqrt{.55(.45)} = 0.4794$$

$$= 0.55 - z_c \frac{(0.4794)}{\sqrt{N}} > 0.5$$

$$0.05 > \frac{z_c (0.4794)}{\sqrt{N}} \Rightarrow \sqrt{N} > \frac{z_c (0.4794)}{0.05} \Rightarrow N > z_c^2 \cdot 99$$

plug in 1.96 $\Rightarrow N > 380.3$
 2.58 $\Rightarrow N > 659$

[This is a bit unrealistic though, because it says $\hat{p} = 0.55$ remains fixed even as we change the sample size N]

(12) If we have binomial form $\sigma_p = \sqrt{\frac{p(1-p)}{N}}$ we can give a more elaborate formula for the CI, which reduces to the usual one for large N.

we know $-z_c \leq \frac{\hat{p} - \mu}{\sigma} \leq z_c$

Rename $\hat{p} \rightarrow Y$
 $\mu \rightarrow p$

$-z_c \leq \frac{Y - p}{\sqrt{\frac{p(1-p)}{N}}} \leq z_c$

$\Rightarrow -z_c \frac{\sqrt{p(1-p)}}{\sqrt{N}} \leq Y - p \leq z_c \frac{\sqrt{p(1-p)}}{\sqrt{N}}$

square everything
 $a < b < c$
 $\Rightarrow a^2 < b^2 < c^2$

$z^2 \frac{p(1-p)}{N} \leq (Y - p)^2 \leq z^2 \frac{p(1-p)}{N}$

so equality: $z^2 \frac{p(1-p)}{N} = Y^2 - 2Yp + p^2$

$\Rightarrow z^2 p - z^2 p^2 = NY^2 - 2NYp + Np^2$

$0 = NY^2 - (z^2 + 2NY)p + (N + z^2)p^2$

Quadratic in p

$p = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{(z^2 + 2NY)}{2(N + z^2)} \pm \frac{\sqrt{(z^2 + 2NY)^2 - 4(N + z^2)NY^2}}{2(N + z^2)}$

$z^4 + 4NYz^2 + 4N^2Y^2 - 4N^2p^2 - 4Np^2z^2$
 $z^2(4NY - 4NY^2 + z^2)$

$= \frac{1}{2N} \frac{(2NY + z^2) \pm z \sqrt{4NY(1-Y) + z^2}}{2N + z^2}$

$= Y + \frac{z^2}{2N} \pm \frac{z}{2N} \sqrt{\frac{4NY}{4N^2} Y(1-Y) + \frac{z^2}{4N^2}}$
 $1 + \frac{z^2}{N}$

Now rename Y back to \hat{p}
p back to μ_p

$\mu_p \approx \frac{\hat{p} + \frac{z^2}{2N} \pm z \sqrt{\frac{p(1-p)}{N} + \frac{z^2}{4N^2}}}{1 + \frac{z^2}{N}}$

$\approx \hat{p} \pm z \sqrt{\frac{p(1-p)}{N}}$
For N large
(they drop $\frac{1}{N}, \frac{1}{N^2}$ terms
but keep $\frac{1}{\sqrt{N}}$)

Does this even have any usefulness?

▷ Confidence Intervals for Sums and Differences

(14) Brand A light bulbs batch size $N_A = 150$
mean lifetime $\bar{X}_A = 1400$
sample std dev $S_A = 120$

Brand B $N_B = 200$
 $\bar{X}_B = 1200$
 $S_B = 80$

Find { 95%
99% } CIs for the difference in lifetimes between the populations

General formula
 $\bar{X} \pm z_c \frac{\sigma}{\sqrt{N}}$

Following ch 8 (see prob 12) $\sigma_{\bar{X}-\bar{Y}} = \sqrt{\frac{\sigma_A^2}{N_A} + \frac{\sigma_B^2}{N_B}}$

so here: $(\bar{X}_A - \bar{X}_B) \pm z_c \sqrt{\frac{S_A^2}{N_A} + \frac{S_B^2}{N_B}} = (1400 - 1200) \pm 1.96 \sqrt{\frac{120^2}{150} + \frac{80^2}{200}}$

Thus we have $\mu_A - \mu_B \in (200 - 22.2, 200 + 22.2)$ with prob 95%.
plug in 2.58 for 99%.

□

(15) A random sample of 400 adults and another sample of 600 teens were asked if they liked a tv show. Construct a CI of 95% for the difference in proportion of the populations of adults and teens.

adults $N_A = 400$
 liked show: $X_A = 100$
 $\hat{p}_A = \frac{100}{400} = \frac{1}{4}$

Teens $N_T = 600$
 liked: $X_T = 300$
 $\hat{p}_T = \frac{300}{600} = \frac{1}{2}$

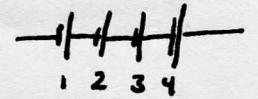
If we do
 $\hat{p}_A - \hat{p}_T = -\frac{1}{4}$

$$(\hat{p}_T - \hat{p}_A) \pm z_c \sqrt{\frac{\hat{p}_A(1-\hat{p}_A)}{N_A} + \frac{\hat{p}_T(1-\hat{p}_T)}{N_T}}$$

From ch 8 results

$$\left(\frac{1}{2} - \frac{1}{4}\right) \pm 1.96 \sqrt{\frac{\frac{1}{4} \cdot \frac{3}{4}}{400} + \frac{\frac{1}{2} \cdot \frac{1}{2}}{600}} \Rightarrow \frac{1}{4} \pm 0.0583 \text{ or } \mu_T - \mu_A \in (0.19, 0.31)$$

(16) EMF of batteries $\mu = 45.1$
 $\sigma = 0.04$

Connect 4 batteries in series 

Let $X_i = \text{emf of battery } i$

Find (a) 95%
 (b) 99%
 (c) 99.73%
 (d) 50% } CIs for total voltage output: $T = \sum_{i=1}^n X_i$

Referring to ch 8 prob #5, we know the first step is to convert to a statement about \bar{x}
 Here we know $T = n\bar{x}$ and we have our one trick $f_{\bar{x}} \approx f_{N(\mu, \sigma^2)}$

All the prev problems have $-z_c \leq \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \leq z_c$

but $\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{\frac{1}{n} \sum X_i - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\frac{1}{n} (\sum X_i - n\mu)}{\frac{\sigma}{\sqrt{n}}} \Rightarrow -z_c \leq \frac{\sum X_i - n\mu}{\sqrt{n}\sigma} \leq z_c$

$$\Rightarrow \underbrace{n\mu}_{4 \cdot 45.1} - z_c \underbrace{\sqrt{n}\sigma}_{2 \cdot 0.04} \leq \sum X_i \leq n\mu + z_c \sqrt{n}\sigma$$

180.4

95% CI $\Rightarrow 180.4 - 1.96(0.08) \leq \sum X_i \leq 180.4 + 1.96(0.08)$

99% CI $\Rightarrow 180.4 \pm 2.58(0.08)$

99.73% $\Rightarrow 180.4 \pm 3(0.08)$

50% $\Rightarrow 180.4 \pm 0.6745(0.08)$

→ This is called 'Probable error'

□