

Introduction to mathematical Statistics

1. Sample size estimation based on the large sample C.I. for p

From the interval $\left[\hat{p} - Z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + Z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$

$$L = \text{length of your } 100(1-\alpha)\% \text{ CI} = 2 \times Z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

L, α, \hat{p} are given and we are interested in sample size n . Therefore,

$$n = \frac{4(Z_{\alpha/2})^2 \hat{p}(1-\hat{p})}{L^2} \leq \frac{4(Z_{\alpha/2})^2 \left(\frac{1}{2}\right)\left(1-\frac{1}{2}\right)}{L^2} = \frac{(Z_{\alpha/2})^2}{L^2}$$

(When $\hat{p} = \frac{1}{2}$, it has the maximum value.)

Example. $L = 0.02, \alpha = 0.05, \hat{p} = 0.54 \Rightarrow n = ?$

Example. $L = 0.02, \alpha = 0.05, \hat{p} = 0.5 \Rightarrow n = ?$

2. Sample size calculation for p based on the maximum error E.

Definition. $P(|\hat{p} - p| \leq E) = 1 - \alpha$

We want to estimate p **within** E with a probability of $(1 - \alpha)$.

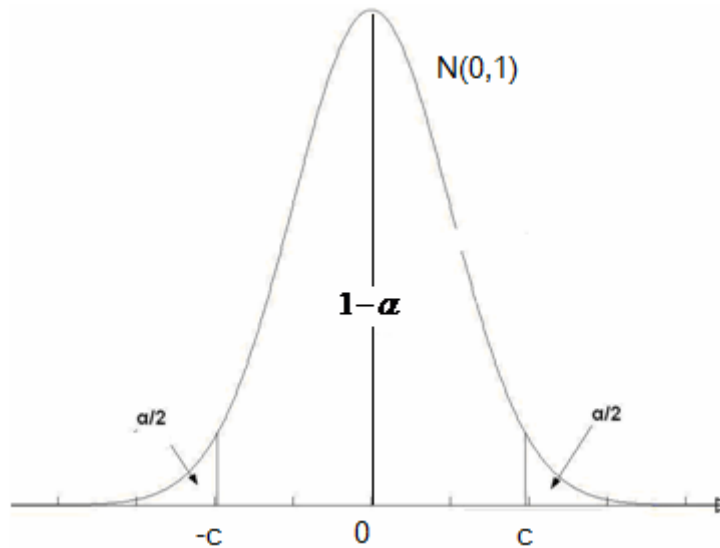
Derive the formula for n

$$P(|\hat{p} - p| \leq E) = 1 - \alpha$$

$$P(-E \leq \hat{p} - p \leq E) = 1 - \alpha$$

$$P\left(-\frac{E}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq \frac{E}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}\right) = 1 - \alpha \text{ and } Z = \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \sim N(0,1)$$

$$\text{Thus: } P\left(-\frac{E}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq Z \leq \frac{E}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}\right) = 1 - \alpha$$



$$c = Z_{\alpha/2} = \frac{E}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}$$

$$\therefore n = \frac{(Z_{\alpha/2})^2 \hat{p}(1-\hat{p})}{E^2} \leq \frac{(Z_{\alpha/2})^2}{4 \cdot E^2}$$

Recall we also derived n based on L - the length of the $100(1-\alpha)\%$ large sample confidence interval for p . Their relationship is

$$L = 2 \cdot E$$

$$P(|\hat{p} - p| \leq E) = 1 - \alpha$$

$$P(-E \leq \hat{p} - p \leq E) = 1 - \alpha$$

$$P(-E - \hat{p} \leq -p \leq E - \hat{p}) = 1 - \alpha$$

$$P(\hat{p} - E \leq p \leq \hat{p} + E) = 1 - \alpha$$

∴ The $100(1 - \alpha)\%$ confidence interval for p is $[\hat{p} - E, \hat{p} + E]$.

The length of the confidence interval is $L = (\hat{p} + E) - (\hat{p} - E) = 2E$

Example. In order to estimate the percent of children with inadequate immunization to be within 0.05 of the true proportion with a probability of 98%

(a) How many children should be sampled?

Solution. $E = 0.05$, $\alpha = 1 - 0.98 = 0.02$

$$\begin{aligned} n &= \frac{(Z_{0.01})^2}{4 \times (0.05)^2} \\ &= \frac{(2.33)^2}{4 \times (0.05)^2} \approx 543 \end{aligned}$$

(b) If the percentage of children with inadequate immunization is estimated to be $\leq 20\%$, then $n = ?$

Solution. $\hat{p} = 20\%$

$$n = \frac{(Z_{0.01})^2 (0.2)(1 - 0.2)}{(0.05)^2} = 348$$

Sample size calculates for 1 population proportions based on the *maximum error E*.

3. Confidence Interval for 1 population mean μ

There are 4 scenarios – we cover only the first 3 scenarios in our class.

1. Normal population, σ^2 is known.

a. Point Estimator : \bar{X}

b. Pivotal Quantity : $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$

c. $100(1-\alpha)\%$ CI for μ : $P(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1-\alpha \Rightarrow \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$

d. Length of CI : $L = 2 \cdot Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

e. Sample size based on L : $n = \frac{4(Z_{\alpha/2})^2 \sigma^2}{L^2}$

f. Sample size based on E

$$P(|\bar{X} - \mu| \leq E) = 1 - \alpha$$

$$P(-E \leq \bar{X} - \mu \leq E) = 1 - \alpha$$

$$P(\bar{X} - E \leq \mu \leq \bar{X} + E) = 1 - \alpha$$

$$\therefore L = (\bar{X} + E) - (\bar{X} - E) = 2E$$

$$\therefore n = \frac{(Z_{\alpha/2})^2 \sigma^2}{E^2}$$

2. Normal population, σ^2 is unknown.

a. Point Estimator : \bar{X}

b. Pivotal Quantity : $T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$

c. $100(1-\alpha)\%$ CI for μ : $P(-t_{n-1,\alpha/2} \leq T \leq t_{n-1,\alpha/2}) = 1-\alpha \Rightarrow \bar{X} \pm t_{n-1,\alpha/2} \cdot \frac{S}{\sqrt{n}}$

d. Length of CI : $L = 2 \cdot t_{n-1,\alpha/2} \frac{S}{\sqrt{n}}$

e. Sample size based on L : $n = \frac{4(t_{n-1,\alpha/2})^2 S^2}{L^2}$

f. Sample size based on E

$$P(|\bar{X} - \mu| \leq E) = 1 - \alpha$$

$$P(-E \leq \bar{X} - \mu \leq E) = 1 - \alpha$$

$$P(\bar{X} - E \leq \mu \leq \bar{X} + E) = 1 - \alpha$$

$$\therefore L = (\bar{X} + E) - (\bar{X} - E) = 2E$$

$$\therefore n = \frac{(t_{n-1,\alpha/2})^2 S^2}{E^2}$$

3. Any populations, large sample

a. Point Estimator : \bar{X}

b. Pivotal Quantity : $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$ or $Z = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim N(0,1)$

c. $100(1 - \alpha)\%$ CI for μ : $P(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1 - \alpha$

$$\Rightarrow \begin{cases} \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \\ \bar{X} \pm Z_{\alpha/2} \cdot \frac{S}{\sqrt{n}} \end{cases}$$

d. Length of CI : $L = 2 \cdot Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$ or $L = 2 \cdot Z_{\alpha/2} \cdot \frac{S}{\sqrt{n}}$

e. Sample size based on L : $n = \frac{4(Z_{\alpha/2})^2 \sigma^2}{L^2}$ or $n = \frac{4(Z_{\alpha/2})^2 S^2}{L^2}$

f. Sample size based on E

$$P(|\bar{X} - \mu| \leq E) = 1 - \alpha$$

$$P(-E \leq \bar{X} - \mu \leq E) = 1 - \alpha$$

$$P(\bar{X} - E \leq \mu \leq \bar{X} + E) = 1 - \alpha$$

$$\therefore L = (\bar{X} + E) - (\bar{X} - E) = 2E$$

$$\therefore n = \frac{(Z_{\alpha/2})^2 \sigma^2}{E^2} \text{ or } n = \frac{(Z_{\alpha/2})^2 S^2}{E^2}$$

4. There also exist other cases, but we don't cover those in our class.

Now I will present more details for Scenario 2 mentioned above.

Scenarios 1 & 3 (easy)

Scenario 2 : normal population, σ^2 unknown

1. Point estimation : $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

2. $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$

3. **Theorem.** Sampling from normal population

- a. $Z \sim N(0,1)$
- b. $W = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$
- c. Z and W are independent.

Definition. $T = \frac{Z}{\sqrt{W/(n-1)}} = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$

----- Derivation of CI, normal population, σ^2 is unknown -----

$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ is not a pivotal quantity.

$\bar{X} - \mu \sim N(0, \frac{\sigma^2}{n})$ is not a pivotal quantity.

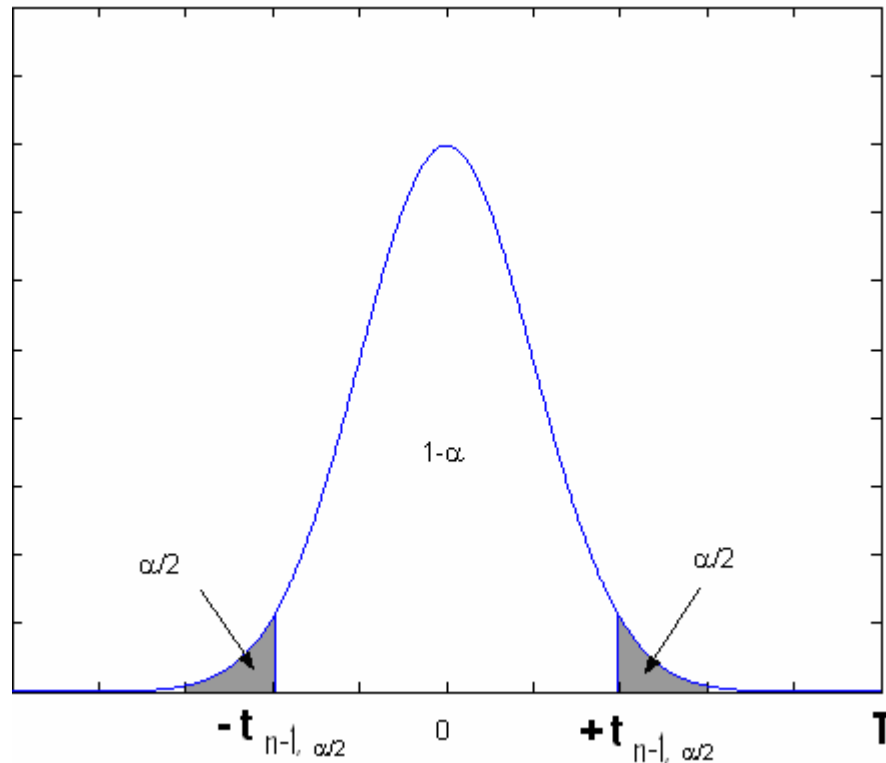
$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$ is not a pivotal quantity.

Remove σ !!!

Therefore $T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$ is a pivotal quantity.

Now we will use this pivotal quantity to derive the $100(1-\alpha)\%$ confidence interval for μ .

We start by plotting the pdf of the t-distribution with $n-1$ degrees of freedom as follows:



The above pdf plot corresponds to the following probability statement:

$$P(-t_{n-1, \alpha/2} \leq T \leq t_{n-1, \alpha/2}) = 1 - \alpha$$

$$\Rightarrow P(-t_{n-1, \alpha/2} \leq \frac{\bar{X} - \mu}{S / \sqrt{n}} \leq t_{n-1, \alpha/2}) = 1 - \alpha$$

$$\Rightarrow P(-t_{n-1, \alpha/2} \frac{S}{\sqrt{n}} \leq \bar{X} - \mu \leq t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}) = 1 - \alpha$$

$$\Rightarrow P(-\bar{X} - t_{n-1, \alpha/2} \frac{S}{\sqrt{n}} \leq -\mu \leq -\bar{X} + t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}) = 1 - \alpha$$

$$\Rightarrow P(\bar{X} + t_{n-1, \alpha/2} \frac{S}{\sqrt{n}} \geq \mu \geq \bar{X} - t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}) = 1 - \alpha$$

$$\Rightarrow P(\bar{X} - t_{n-1, \alpha/2} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}) = 1 - \alpha$$

\Rightarrow Thus the $100(1 - \alpha)\%$ C.I. for μ when σ^2 is unknown is

$$\left[\bar{X} - t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}, \bar{X} + t_{n-1, \alpha/2} \frac{S}{\sqrt{n}} \right]. \quad (*\text{Please note that } t_{n-1, \alpha/2} \geq Z_{\alpha/2} *)$$

Example. In a random sample of $n = 36$ parochial schools throughout the south, the average number of pupils per school is 379.2 with a standard deviation of 124. Use the sample to construct a 95% CI for μ , the mean number of pupils per school for all parochial schools in the south.

Solution. CI for μ , large sample

$$n = 36, \bar{X} = 379.2, S = 124, \alpha = 0.05$$

$$95\% \text{ CI for } \mu \text{ is } \bar{X} \pm Z_{0.05} \cdot \frac{S}{\sqrt{n}} = 379.2 \pm 1.96 \cdot \frac{124}{\sqrt{36}}$$

$$\therefore [338.7, 419.7]$$

Example. In a psychological **depth-perception** test, a random sample of $n = 14$ airline pilots were asked to judge the distance between 2 markers at the other end of a laboratory.

The data (in test) are

$$2.7, 2.4, 1.9, 2.4, 1.9, 2.3, 2.2, 2.5, 2.3, 1.8, 2.5, 2.0, 2.2, 2.6$$

Please construct a 95% CI for μ , the average distance.

Solution.

(Note: we can perform the Shapiro-Wilk test to examine whether the sample comes from a normal population or not. This test is not required in our class. Here we simply assume the population is normal. I will always give you such information in the exams.)

CI for μ , small sample, normal population, population variance unknown.

$$n = 14, \bar{X} = 2.26, S = 0.28, \alpha = 0.05$$

$$95\% \text{ CI for } \mu \text{ is } \bar{X} \pm t_{n-1, \alpha/2} \cdot \frac{S}{\sqrt{n}} = 2.26 \pm 2.16 \cdot \frac{0.28}{\sqrt{14}}$$

$$\therefore [2.10, 2.42]$$

Example. A federal agency has decided to investigate the advertised **weight** we printed on cartons of a certain brand of cereal. Historical data show that $\sigma = 0.75$ ounce. If we wish to estimate the weight within 0.25 ounce with 99% confidence, how many cartons should be sampled?

Solution. $E = 0.25$, $\sigma = 0.75$, $\alpha = 0.01$

$$n = \frac{(Z_{0.005})^2 (0.75)^2}{(0.25)^2} \approx 60$$