# <u>Confidence Interval</u>

### *<u>Illustrated through inference on one population mean or proportion</u>*

## Motivation & simple random sample

Eg) We wish to estimate the average height of adult US males

➔ Take a random sample.

- "Simple" random sample: every subject in the population has the same chance to be selected.

## Introduction to statistical inference on one population mean

For a "<span style="color:blue">random</span> sample" of size n: $X_1, X_2, \ldots, X_n$

**<i> Point estimatior $\overline{X}$ → sample mean** $(= \dfrac{X_1 + X_2 + \ldots + X_n}{n} = \dfrac{\sum\limits_{i=1}^{n} X_i}{n})$

Other estimators: median, mode, trimmed mean, …

**<ii> Confidence Interval (C.I.)**

Eg)  95% C.I. for μ

99.9999% C.I. ('6-9' in the manufacture industry)

**<iii> Hypothesis Test**

Eg) $H_0 : \mu \le 5'6''$

$H_1 : \mu \succ 5'6''$

Point Estimator, C.I., Test $\Rightarrow$ Statistical Inference

- Draw some conclusion on the population (parameters of interest) based on a random sample.

# 1. The Exact Confidence Interval for <mark>μ</mark> when the population is normal & σ² is known

---

① **Point estimator** and **confidence interval** for $\mu$

- When the population is normal and the population variance is known.
- Let $X_1, X_2, \ldots, X_n$ be a random sample for a normal population with mean $\mu$ and variance $\sigma^2$. That is, $X_i \overset{iid.}{\sim} N(\mu, \sigma^2), i = 1, \ldots, n$.
- For now, we assume that $\sigma^2$ is known.

---

**<i> Point Estimator for $\mu$ :** $\hat{\mu} = \overline{X} \sim N(\mu, \dfrac{\sigma^2}{n})$

$E(\hat{\mu}) = E(\overline{X}) = \mu \Rightarrow \hat{\mu} = \overline{X}$ is an unbiased estimator of $\mu$

- Intuitively, this means if you take "many" samples of size n from the population, then the mean of these samples means would be equal to $\mu$ if you take a large enough # of samples.

- $\overline{X}$ is also a maximum likelihood estimator (MLE) of $\mu$.

- $\overline{X}$ is also a method of moment estimator (MOME) of $\mu$.

- Other good properties too.


**<ii> Confidence Interval for $\mu$**

- Intuitive approach (backwards derivation for the CI boundaries $C_1 \, and \, C_2$):

$$P(C_1 \leq \mu \leq C_2) = 0.95$$

$$P(-C_1 \geq -\mu \geq -C_2) = 0.95$$

$$P(\overline{X} - C_1 \geq \overline{X} - \mu \geq \overline{X} - C_2) = 0.95$$

$$P\left(\frac{\overline{X} - C_1}{\sigma/\sqrt{n}} \geq \frac{\overline{X} - \mu}{\sigma/\sqrt{n}} \geq \frac{\overline{X} - C_2}{\sigma/\sqrt{n}}\right) = 0.95$$

*Since we know*

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

We can compute the expressions for $C_1$ and $C_2$.

However, one question is that there are MANY ways to choose the C's.

Later you will see that for pivotal quantity with symmetric pdfs, the symmetric CIs are the optimal – in that they have the shortest lengths for the given confidence level $100(1-\alpha)\%$.

## Now we present a general approach to derive the CI's.

**General approach for deriving CI's : the Pivotal Quantity (P.Q.) approach**

*Definition: A pivotal quantity is a function of the sample and the parameter of interest. Furthermore, its distribution is entirely known.

1. We start by looking at the point estimator of $\mu$. $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

* Is $\bar{X}$ a pivotal quantity for $\mu$?

$\rightarrow \bar{X}$ is not because $\mu$ is unknown.

* function of $\bar{X}$ and $\mu$ : $\bar{X} - \mu \sim N(0, \frac{\sigma^2}{n})$

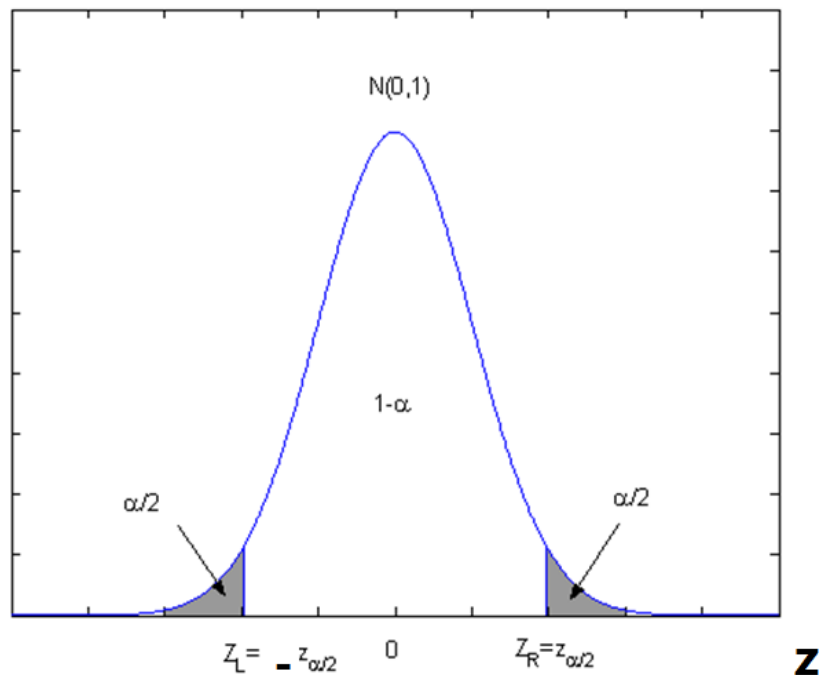$\rightarrow$ Yes, it is pivotal quantity.

* Another function of $\bar{X}$ and $\mu$ : $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$

$\rightarrow$ Yes, it is pivotal quantity.

So, Pivotal Quantity is not unique.

2. Now that we have found the pivotal quantity Z, we shall start the derivation for the symmetrical CI's for μ from the PDF of the pivotal quantity Z

N(0,1)

$1-\alpha$

$\alpha/2$

$\alpha/2$

$Z_L = -z_{\alpha/2}$   0   $Z_R = z_{\alpha/2}$   **Z**

$100(1-\alpha)\%$ CI for $\mu$, $0<\alpha<1$

(e.g. $\alpha=0.05 \Rightarrow 95\%$ C.I.)

$$P(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1-\alpha$$

$$P(-Z_{\alpha/2} \leq \frac{\overline{X}-\mu}{\sigma/\sqrt{n}} \leq Z_{\alpha/2}) = 1-\alpha$$

$$P(-Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \overline{X} - \mu \leq Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1-\alpha$$

$$P(-\overline{X} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq -\mu \leq -\overline{X} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1-\alpha$$

$$P(\overline{X} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \geq \mu \geq \overline{X} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1-\alpha$$

$$P(\overline{X} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \overline{X} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1-\alpha$$

$\therefore$ the $100(1-\alpha)\%$ C.I. for $\mu$ is $[\overline{X} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \overline{X} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}]$

\*Note, some special values for $\alpha$ and the corresponding $Z_{\alpha/2}$ values are:

1. The 95% CI, where $\alpha = 0.05$ and the corresponding $Z_{\frac{\alpha}{2}} = Z_{0.025} = 1.96$

**Example 1.** A random sample of 400 adult US male was taken and the sample mean was found to be $\bar{X} = 5'7'' = 67 \ inches$ . Based on past studies, it is believed that the population distribution of all adult US male is normal and the standard deviation is 30 inches. Please construct a 95% confidence interval for the average height of all adult US male based on this sample.

    **Solution:** The 95% CI for $\mu$ is

$$\left[\bar{X} - 1.96\frac{\sigma}{\sqrt{n}}, \bar{X} + 1.96\frac{\sigma}{\sqrt{n}}\right] = \left[67 - 1.96\frac{30}{\sqrt{400}}, 67 + 1.96\frac{30}{\sqrt{400}}\right] \approx [64, 70]$$

    That is, the estimated 95% confidence interval for the average height of all adult US male is $[5'4'', \ 5'10'']$.
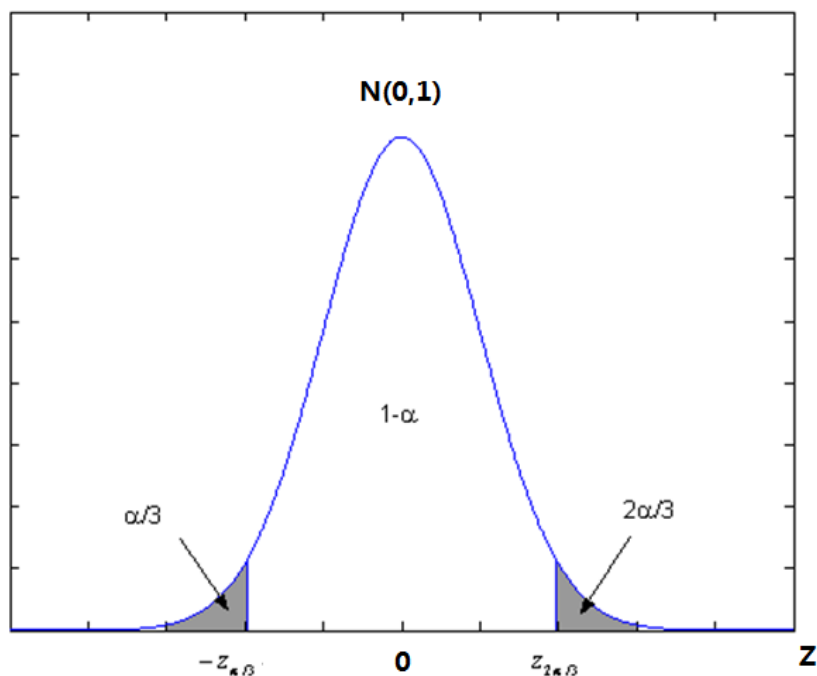
        ...                      ...

This means that we are 95% sure the population mean $\mu$ would lie between $5'4''$ and $5'10''$.

$\therefore$ Recall the 100(1-$\alpha$)% symmetric C.I. for $\mu$ is $[\overline{X} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \overline{X} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}]$

**\*Please note that this CI is symmetric around $\overline{X}$**

**The length of this CI is:** $L_{sy} = 2 \cdot Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$

Now we derive a <u>non-symmetrical CI</u>:

**N(0,1)**

$1-\alpha$

$\alpha/3$

$2\alpha/3$

$-Z_{\alpha/3}$ ∙ **0** $Z_{2\alpha/3}$ **Z**

$P(-Z_{\alpha/3} \le Z \le Z_{2/3\alpha}) = 1 - \alpha$

100(1-α)% C.I. for μ

$$\Rightarrow [\overline{X} - Z_{2/3\alpha} \cdot \frac{\sigma}{\sqrt{n}}, \overline{X} + Z_{1/3\alpha} \cdot \frac{\sigma}{\sqrt{n}}]$$

---

Compare the lengths of the C.I.'s, one can prove theoretically that:

$$L = (Z_{\alpha/3} + Z_{2/3\alpha}) \cdot \frac{\sigma}{\sqrt{n}} > L_{sy} = 2 \cdot Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

**You can try a few numerical values for α, and see for yourself. For example,**

**α = 0.05**

---

HW: Please derive the 100(1-α)% symmetric C.I. for μ based on a random sample from a normal population with unknown variance

## 2. (Large Sample) Confidence interval for a population mean (*any population) or a population proportion p

**<Theorem> Central Limit Theorem**

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \xrightarrow{n \to \infty} N(0,1)$$

**When n is large enough, we have**

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

That means Z follows approximately the normal (0,1) distribution.

**Application #1. Inference on $\mu$ when the population distribution is unknown but the sample size is large**

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

By Slutsky's Theorem We can also obtain another pivotal quantity when σ is unknown by plugging the sample standard deviation S as follows:

$$Z = \frac{\bar{X} - \mu}{S/\sqrt{n}} \overset{\cdot}{\sim} N(0,1)$$

We subsequently obtain the $100(1-\alpha)\%$ C.I. using the second P.Q. for $\mu$ : $\quad \bar{X} \pm Z_{\alpha/2} \dfrac{S}{\sqrt{n}}$

**Application #2. Inference on one population proportion p when the population is Bernoulli(p) *** ** Let $X_i \overset{i.i.d.}{\sim} Bernoulli(p), \ i = 1, \cdots, n$ , please find the $100(1-\alpha)\%$ CI for p.

Point estimator : $\hat{p} = \bar{X} = \dfrac{\sum_{i=1}^{n} X_i}{n}$ (ex. $n = 1000$ , $\hat{p} = 0.6$ )

Our goal: derive a $100(1-\alpha)\%$ C.I. for p

Thus for the Bernoulli population, we have:

$$\mu = E(X) = p$$

$$\sigma^2 = Var(X) = p(1-p)$$

*Thus by the CLT we have:*

$$Z = \frac{\bar{X} - p}{\sqrt{\frac{p(1-p)}{n}}} \overset{\cdot}{\sim} N(0,1)$$

Furthermore, we have for this situation: $\bar{X} = \hat{p}$

Therefore we obtain the following pivotal quantity Z for p:

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \overset{\cdot}{\sim} N(0,1)$$

By Slustky's theorem, we can replace the population proportion in the denominator with the sample proportion and obtain another pivotal quantity for p:

$$Z^* = \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \overset{\cdot}{\sim} N(0,1)$$

**# Thus the $100(1-\alpha)$% (approximate, or large sample) C.I. for p based on the second pivotal quantity $Z^*$ is:**

$$P(-z_{\alpha/2} \leq Z^* \leq z_{\alpha/2}) = 1 - \alpha$$

$$P\left(-z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

$$P\left(-\hat{p} - z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq -p \leq -\hat{p} + z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 1 - \alpha$$

$$P\left(\hat{p} - z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 1 - \alpha$$

=> The $100(1-\alpha)$% large sample C.I. for p is

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

# CLT => n large usually means $n \geq 30$

# special case for the inference on p based on a Bernoulli population. The sample size n is large means

Let $X = \sum_{i=1}^{n} X_i$, large sample means:

$n\hat{p} = X \geq 5$ (*Here X= total # of 'S'), and

$n(1-\hat{p}) = n - X \geq 5$ (*Here n-X= total # of 'F')

## Example 2.

During one of the "beer wars" in the early 1980's, a taste test between Schlitz and Budweiser was the focus of a TV commercial. 100 people agreed to drink 2 unmarked mugs and indicate which of the two beers they liked better. 54 chose "Bud". Construct and interpret the corresponding 95% confidence interval for $p$ - the proportion of beer drinkers who prefer Bud to Schlitz.

## Solution.

**Confidence Interval for one population proportion (p) when the sample size is large**

Sample size : $n$ $(n = 100)$

Sample proportion : $\hat{p} = \dfrac{\sum_{i=1}^{n} X_i}{n}$ ( $\hat{p} = \dfrac{54}{100}$ )

*** Recall we usually denote $X = \sum_{i=1}^{n} X_i$

**"sample is large"** means

- For one population mean, $n \geq 30$

- For one population proportion : $X \geq 5$ and $(n - X) \geq 5$

$$\left( X = 54 \geq 5 \, ; \, n - X = 46 \geq 5 \right)$$

$n = 100, \, X = 54, \,$ 95% CI for $p$

From 95% confidence interval, $1 - \alpha = 0.95, \, \alpha = 0.05, \, \dfrac{\alpha}{2} = 0.025$

$$\hat{p} = \frac{54}{100} = 0.54 \, ; \, Z_{0.025} = 1.96$$

$$\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{(0.54)(0.46)}{100}} = 0.049$$

$$Z_{0.025} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 1.96 \times 0.049 = 0.096$$

$\therefore$ The 95% confidence interval for $p$ is $\left[ 0.444, 0.636 \right]$

If $n = 10000 \, ; \, \hat{p} = 0.54$ ,

$$\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{(0.54)(0.46)}{10,000}} = 0.0049$$

$$Z_{0.025} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 1.96 \times 0.0049 = 0.0096 \approx 0.01$$

$\therefore$ The 95% confidence interval for $p$ is $\left[ 0.53, 0.55 \right]$

# 3. The Exact Confidence Interval for μ when the population is normal & σ² is unknown

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

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

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

   a. $Z \sim N(0,1)$

   b. $W = \dfrac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$

   c. $Z$ and $W$ are independent.

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

------ **Derivation of CI, normal population, $\sigma^2$ is unknown** ------

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

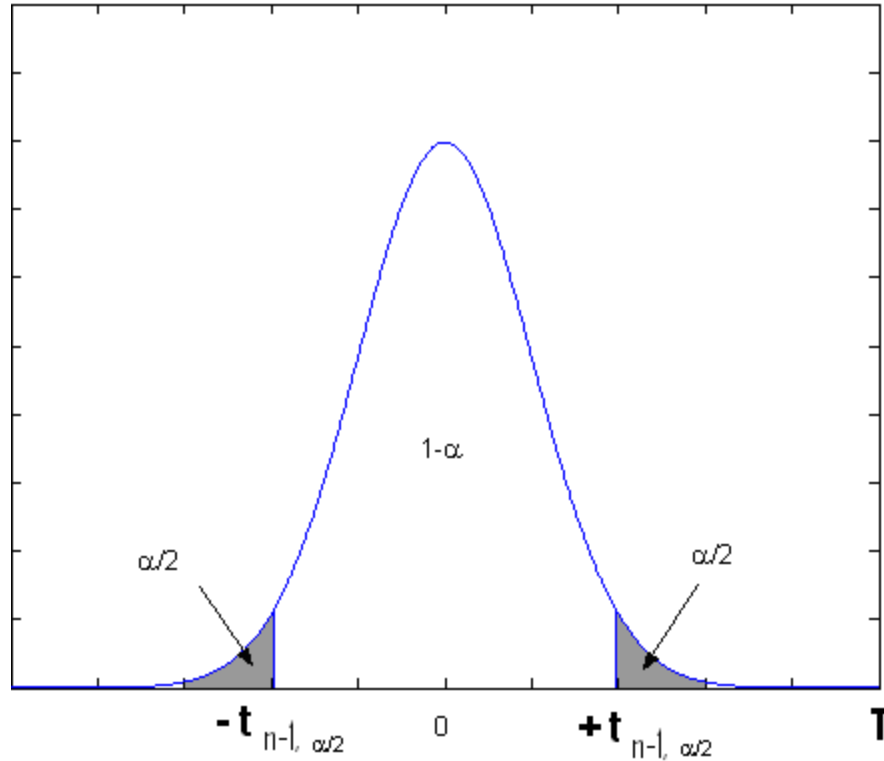$\overline{X} - \mu \sim N(0, \dfrac{\sigma^2}{n})$ is not a pivotal quantity.

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

Remove $\sigma$ !!!

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

Now we will use this pivotal quantity to derive the 100(1-α)% 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:

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

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

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

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

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

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

=> Thus the $100(1-\alpha)\%$ C.I. for $\mu$ when $\sigma^2$ is unknown is

$$[\bar{X} - t_{n-1,\alpha/2} \frac{S}{\sqrt{n}}, \bar{X} + t_{n-1,\alpha/2} \frac{S}{\sqrt{n}}].$$

(*Please note that $t_{n-1,\alpha/2} \geq Z_{\alpha/2}$)

**Example 3.** 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 $\mu$, 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 $\mu$, small sample, normal population, population variance unknown.

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

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

$\therefore \left[ 2.10, 2.42 \right]$