# Confidence Intervals

## Some ~~parameter~~ Terms

1. **Population**

2. **Sample**

3. **Parameter :-** Population $\rightarrow$ mean($\mu$) $\rightarrow$ var $\Big\}$ parameter $\rightarrow$ Std($\sigma$)
   $\uparrow$ sigma

   Since it is often difficult or impossible to obtain data from an entire population, parameters are usually unknown and must be estimated based on available sample data.

4. **Statistic ?-** A statistic is a numerical value that describe a characteristic of a sample, which is a subset of the population. By using statistic calculated from a representative sample, researchers can make inferences about the unknown respective parameter of the population. common statistic include the sample mean (denoted by $\bar{x}$, pronounced "$x$-bar"), the sample median, and the sample standard deviation (denoted by $s$).

5. **Inferential Statistics**

   Sample $\longrightarrow$ predict ~~prop of~~ about population.

   Inferential statistics use hypothesis testing, confidence interval and regression analysis, among others.

   These methods help researchers answer question like:
   a) is there a significant difference btn two group?
   b) can we predict the outcome of a variable based on the value of other values?

c. What is the relationship b/w two or more variable?

## Point Estimate

Subs → (77 k) → any age?

or

Live → (100) → $\bar{X}$ → Sample mean (point estimate)
velmer

10 blue classes
→ (100) → 10 mean ($\bar{X}_1, \bar{X}_2, \bar{X}_3 \cdots \bar{X}_{10}$)
velmer

$\boxed{\bar{X}_1, \bar{X}_2, \bar{X}_3 \cdots \bar{X}_{10}}$ → $\bar{X}_{mean}$ (point estimate)

A point estimate is a single value calculated from a sample, that serves as the best guess or approximation for an unknown population parameter, such as the mean or standard deviation. Point estimates are often used in statistics when we want to make inference about a population based on a sample.

# Confidence Interval

**Confidence Interval:** In simple words, is a range of values within which we expect a particular <u>population parameter</u>, like a mean, to fall. It's a way to express the uncertainty around an estimate obtained from a sample of data.

$$\mu, \sigma$$

**Confidence level :** usually expressed as a percentage like 95% Indicates how sure we are that the true value lies within the interval.

$$[25, 32]$$
$\longrightarrow$ Confidence Interval

95% of data (confidence level)
both confidence Interval.

Confidence Interval = $[Point\ Estimate] \pm [Margin\ of\ error]$

$$25 \quad \pm \quad 4$$

25+4 = 29
25-4 = 21

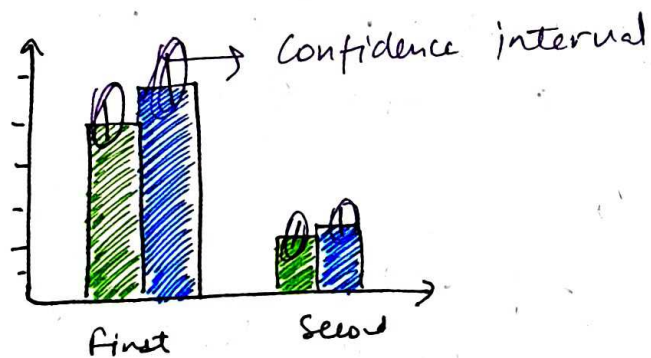Confidence Interval = $[21, 29]$

**Ways to calculate CI:**

z procedure
pop $\rightarrow$ std (available)
$(\sigma)$

t procedure
pop $\rightarrow$ std (not available)
$(\sigma)$

* Confidence interval is created for <u>parameter</u> and not statistics. Statistics help us get the confidence interval for a parameter

## Example of CI usage


→ Confidence interval

First    Second

## Confidence Interval (Sigma known)
↳ Z procedure

## Assumption
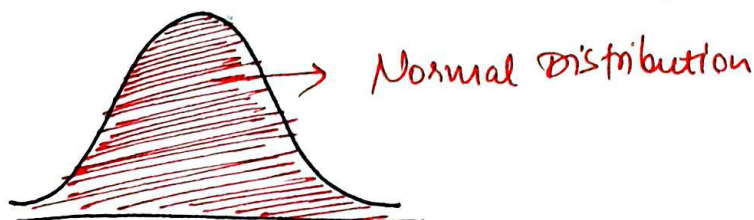
1. <u>Random Sampling</u>

   pop ⟶ <u>Sample</u>
   ↳ sample number random.

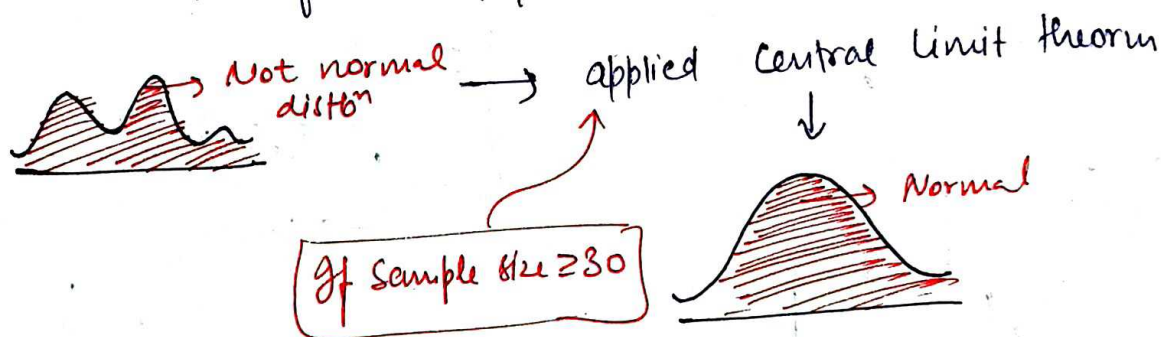   e.g → Collect salary data of all India. So, we have to collect data from Random state <u>not only</u> <u>single state.</u>

2. Known population standard deviation

3. <u>Normal distribution</u> or large sample size


→ Normal Distribution

However, if the population distⁿ is not normal, the Central limit theorm can be applied when the sample size is large ( usually, sample size $n \geq 30$ is considered large enough). According to the Central limit Theorm, the sampling distⁿ of the sample mean will approach a normal distⁿ as the sample size increase, regardless of the shape of the population distⁿ.



A $(1 - alpha) * 100\%$ confidence interval for mu:

YT $\rightarrow$ Campus% $\rightarrow$ 77K $\rightarrow$ 28 ± 14

[16, 42] ← Confidence interval

$\sigma = 15$

↳ 95% confidence level

formula CI using Z procedure

$$\text{Point estimate } (\overline{X}) \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$(1 - \alpha)$ ⟶ Confidence level (95%)

$\sigma$ ⟶ std pop

$n$ ⟶ Sample size ⟶ 100
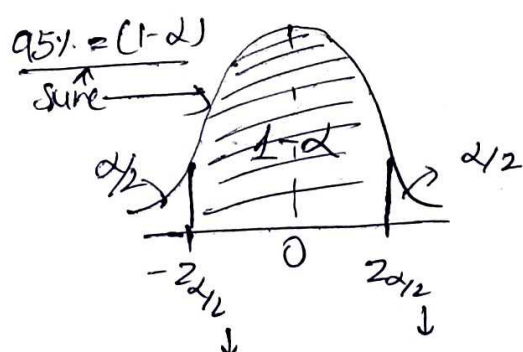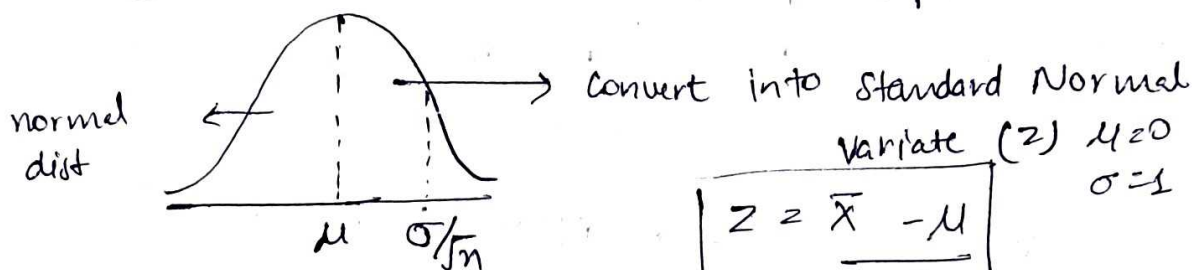
# Intuition

Point estimate ($\bar{x}$) → CLT

10 live class → (50 Student) → Avg age

$$\boxed{\bar{X}_{1\,class}, \bar{X}_{2\,class}, - - - - - \bar{X}_{10\,class}}$$ → Sampling dist of Sample mean.

normal dist

Convert into standard Normal variate (Z) $\mu = 0$, $\sigma = 1$

$$\boxed{Z = \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}}}$$

└→ graph

95% = (1-$\alpha$) Sure →

$$1 - \alpha$$

$\alpha/2$

$\alpha/2$

$-Z_{\alpha/2}$

$Z_{\alpha/2}$

$1 - \alpha = 95\%$   $-Z_{2.5} \leftarrow \boxed{\phantom{}} \rightarrow Z_{2.5}$

$\boxed{\alpha = 5}$   95%.

$\Rightarrow$

$$CI = \bar{X} \pm \underset{95\%}{\left( Z_{\alpha/2} \right)} \dfrac{\sigma}{\sqrt{n}}$$

$$P\left( -Z_{\alpha/2} < Z < Z_{\alpha/2} \right) = 1 - \alpha$$

$$P\left( -Z_{\alpha/2} < \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha/2} \right) = 1 - \alpha$$

$$P\left( -Z_{\alpha/2}\dfrac{\sigma}{\sqrt{n}} < \bar{X} - \mu < Z_{\alpha/2}\dfrac{\sigma}{\sqrt{n}} \right) = 1 - \alpha$$

$$P\left( -\bar{X} - Z_{\alpha/2}\dfrac{\sigma}{\sqrt{n}} < -\mu < -\bar{X} + Z_{\alpha/2}\dfrac{\sigma}{\sqrt{n}} \right) = 1 - \alpha$$

$$P\left( \bar{X} - Z_{\alpha/2}\dfrac{\sigma}{\sqrt{n}} < \mu < \bar{X} + Z_{\alpha/2}\dfrac{\sigma}{\sqrt{n}} \right) = 1 - \alpha$$

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

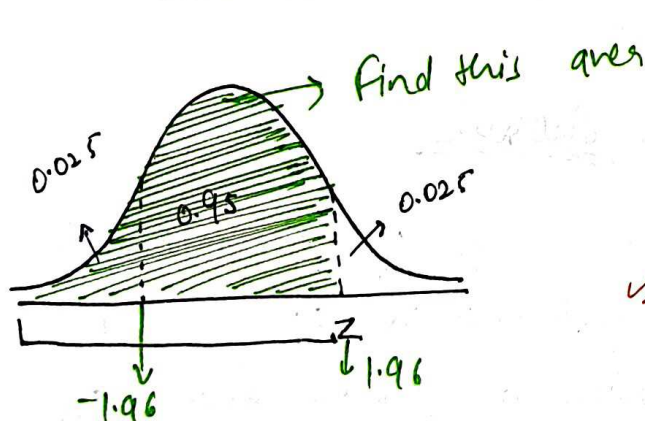$\overline{x} \to$ Sample

* $\mu$ is a fixed (Age is fixed) so, prob is also fixed.
($\overline{x}$) is sample and it is not fixed bcz
every time ($\overline{x}$) sample always change. We cannot
say " sample $\mu$ lie both $\overline{x}-Z_{\alpha/2}\frac{\sigma}{\sqrt{n}}$ and $\overline{x}+Z_{\alpha/2}\frac{\sigma}{\sqrt{n}}$ " prob
is 95%".

we can say " $\mu$ lie between $\overline{x}-Z_{\alpha/2}\frac{\sigma}{\sqrt{n}}$ and $\overline{x}+Z_{\alpha/2}\frac{\sigma}{\sqrt{n}}$
confidence level is 95%".

CI $\quad \mu = \overline{x} \pm Z_{\alpha/2}\frac{\sigma}{\sqrt{n}}$ $\qquad (1-\alpha) \to 95\%$.

$\qquad \alpha/2 = \frac{0.50}{2}$

$\qquad \alpha/2 = 0.250$

$$\boxed{\mu = \overline{X} \pm Z_{0.250}\frac{\sigma}{\sqrt{n}}}$$

$\to$ find this area $\qquad$ $\frac{950}{025}$ $\underbrace{(975)} \to$ find this value

on Z table.



0.025 $\quad$ 0.95 $\quad$ 0.025

$-1.96 \qquad Z$
$\qquad\qquad 1.96$

$$\boxed{\mu = \overline{X} \pm 1.96\frac{\sigma}{\sqrt{n}}}$$

$\boxed{Z_{\alpha/2} = 1.96}$

* If I want to to $Z_{\alpha/2}$ of 75% then what we have
to do is.

$\qquad\qquad\qquad$ Area using Z score.



0.125 $\quad$ 0.75 $\quad$ 0.125

**\* What is confidence Interval?**

→ Confidence Interval is not a Probability to defined interval [14-42] or not a ~~any~~ any lie both there [14-42].

<u>real defination</u>

↳ 95% ⟷ [18, 42]
      ↓
  Channel

⟨77k⟩ → 100 times ( Random Samples )(50)

avg age of 50 people →      1 time   ⎫  100 time mai
                               ⎬ → se 95 time
                                 ⎭   of avg age

avg age of 50 people →      100times    [18, 42) ke
                                         beech mai ayega

<u>Interpreting Confidence Interval</u>

A confidence interval is a range of values within which a population parameter, such as the pop mean, is estimated to lie with a certain level of confidence. The confidence interval provides an indication of the precision and uncertainty associated with the estimate. To interpret the confidence interval values, consider the following points : -

1. Confidence ~~Interval~~ Level: The confidence level (commonly set at 90%, 95% or 99%) represents the prob that the confidence interval will contain the true pop parameter if the sampling and estimation process were repeated multiple times. For example, a 95% confidence interval means that if you were to draw 100 different samples from the pop and calculate the confidence interval for each, approximately 95 of those interval would contain the true pop parameter.

2. Interval range: The width of the confidence interval gives an indication of the precision of the estimation. A narrower confidence interval suggests a more precise estimate of the pop parameter, while a wider interval indicates greater uncertainty. The width of the interval depends on the sample size, variability in the data, and the desired level of confidence.

3. Interpretation: To interpret the confidence interval values, you can say that you are "XX confident that the true pop parameter lies within the range (lower limit, upper limit)". Keep in mind that this statement is about the interval, not the specific point estimate, and it refers to the confidence level you choose when constructing the interval.

# Factor Affecting Margin Error

1. Confidence level $(1-alpha)$.   $\dfrac{Upper - Lower}{2} = $ margin of error.

CI = point estimate $\pm$ margin of error

$$= \bar{X} \pm Z_{\alpha/2} \dfrac{\sigma}{\sqrt{n}}$$

Sample mean ← $\bar{X}$

Critical ① Value → $Z_{\alpha/2}$

Pop std ② → $\sigma$

Sample size ③ → $\sqrt{n}$

① Confidence level ↓   Margin of error ↓

en →   50% → [49.09, 50]

95% → [39.01, 55.01]

② Standard deviation ↓   Margin of error ↓

③ Sample size ↑   Margin of error ↓

# Confidence Interval (Sigma not Known)

Using the t procedure

## Assumption

1. **Random Sampling:** The data must be collected using a random sampling method to ensure that the sample is representative of the pop. This helps to minimize biases and ensures that the results can be generalized to the entire pop.

2. **Sample Standard deviation:** The pop standard deviation ($\sigma$) is unknown, and the sample standard deviation (S) is an estimate. The t-dist is ~~specially~~ Specifically designed to account for the additional uncertainty introduced by using the sample standard deviation instead of the pop standard deviation.

3. **Approximately normal Distribution:** The t-procedure assumes that the underlying pop is approximately normally distributed, or the sample size is larger enough for the central limit Theorem to apply. If the pop distⁿ is heavey skewed or has extreme outliers, the t-procedure may not be accurate and non-parametric methods should be considered.

4. **Independent Observation:** The Observations in the sample should independent of each other. In other words, the value of one observation should not influence the value of another observation. This is particularly imp when working with time
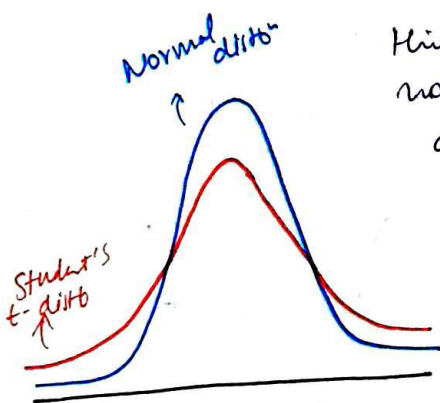
series data or data with inherent dependencies.

In CIL with signal $= \bar{X} \pm Z_{\alpha/2} \underbrace{\sigma}_{\sqrt{n}} \xrightarrow{\times} s$

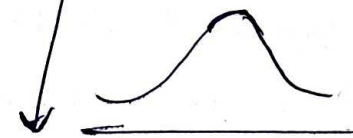$$CI = \bar{X} + Z_{\alpha/2} \frac{\textcircled{s}}{\sqrt{n}} \rightarrow \text{create complexity}$$

$\bar{X} \rightarrow \boxed{\dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}}} \rightarrow$ Std normal dist

$Z' \rightarrow \boxed{\dfrac{\bar{X} - \mu}{S/\sqrt{n}}} \longrightarrow$ Student's t distribution   ( look like normal disth but not normal distn)

This is not normal disth

Normal disth



Student's t distrb

theoritical disth (Not exist in nature like Normal disth)

Only one parameter → degree of freedom $(n-1)$

degree of freedom $- \infty$
   ↳ Student's t disth = Normal distribution

degree of freedom ↑↑ (Student's t disth ≈ Normal disth)

$$CI = \bar{X} + \textcircled{t}_{\alpha/2} \frac{S}{\sqrt{n}}$$

use t table

95% of $Z_{\alpha/2} \rightarrow 1.96$

95% of $t_{\alpha/2} \approx 2.0008$

always $t_{\alpha/2} > Z_{\alpha/2}$ bcz we are not confident so increase the interval for more confidence.