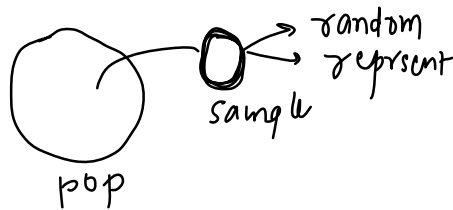# Some Terms

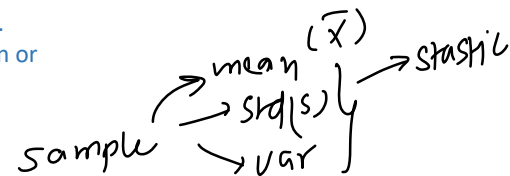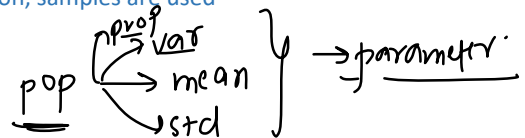**Population Vs Sample**

Population: A population is the entire group or set of individuals, objects, or events that a researcher wants to study or draw conclusions about. It can be people, animals, plants, or even inanimate objects, depending on the context of the study. The population usually represents the complete set of possible data points or observations.

Sample: A sample is a subset of the population that is selected for study. It is a smaller group that is intended to be representative of the larger population. Researchers collect data from the sample and use it to make inferences about the population as a whole. Since it is often impractical or impossible to collect data from every member of a population, samples are used as an efficient and cost-effective way to gather information.
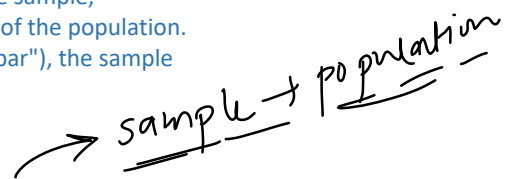
**Parameter Vs ~~Estimate~~**

Parameter: A parameter is a numerical value that describes a characteristic of a population. Parameters are usually denoted using Greek letters, such as μ (mu) for the population mean or σ (sigma) for the population standard deviation. 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.

$$\bar{X} \rightarrow \mu$$

Statistic: A statistic is a numerical value that describes a characteristic of a sample, which is a subset of the population. By using statistics calculated from a representative sample, researchers can make inferences about the unknown respective parameter of the population. Common statistics include the sample mean (denoted by x̄, pronounced "x-bar"), the sample median, and the sample standard deviation (denoted by s).

**Inferential Statistics**

Inferential statistics is a branch of statistics that focuses on making predictions, estimations, or generalizations about a larger population based on a sample of data taken from that population. It involves the use of probability theory to make inferences and draw conclusions about the characteristics of a population by analysing a smaller subset or sample.

The key idea behind inferential statistics is that it is often impractical or impossible to collect data from every member of a population, so instead, we use a representative sample to make inferences about the entire group. Inferential statistical techniques include hypothesis testing, confidence intervals, and regression analysis, among others.

These methods help researchers answer questions like:

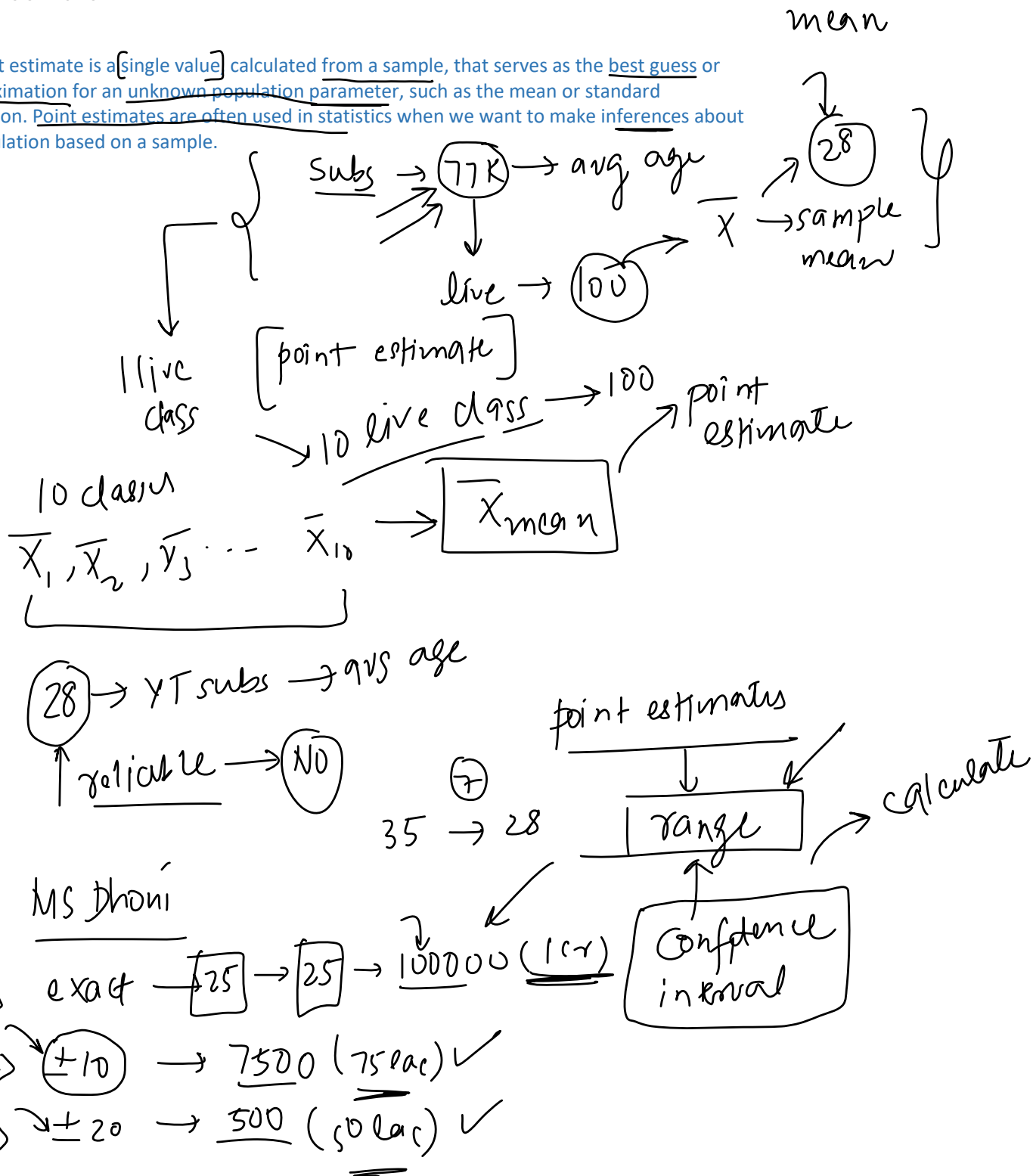    a. Is there a significant difference between two groups?
    b. Can we predict the outcome of a variable based on the values of other variables?
    c. What is the relationship between two or more variables?

Inferential statistics are widely used in various fields, such as economics, social sciences, medicine, and natural sciences, to make informed decisions and guide policy based on limited data.

# 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 inferences about a population based on a sample.

mean

Subs → (77K) → avg age

live → (100) → $\overline{X}$ → sample mean

(28)

1 live class

[point estimate]

→ 10 live class → 100 → point estimate

10 classes

$\overline{X}_1, \overline{X}_2, \overline{Y}_3 \cdots \overline{X}_{10}$ → $\boxed{\overline{X}_{mean}}$

(28) → YT subs → avg age

reliable → (NO)

point estimates

35 → 28   $\boxed{range}$ → calculate

Confidence interval

MS Dhoni

1) exact → $\boxed{25}$ → $\boxed{25}$ → 100000 (1cr)

2) ±10 → 7500 (75 lac) ✓

3) ±20 → 500 (50 lac) ✓

# Confidence Interval

30 March 2023    07:18

$\mu, \sigma$    $\overline{X}$    $\dfrac{2}{2}$    $[25, 32]$

**Confidence interval**, in simple words, is a range of values within which we expect a particular population parameter, like a mean, to fall. It's a way to express the uncertainty around an estimate obtained from a sample of data.
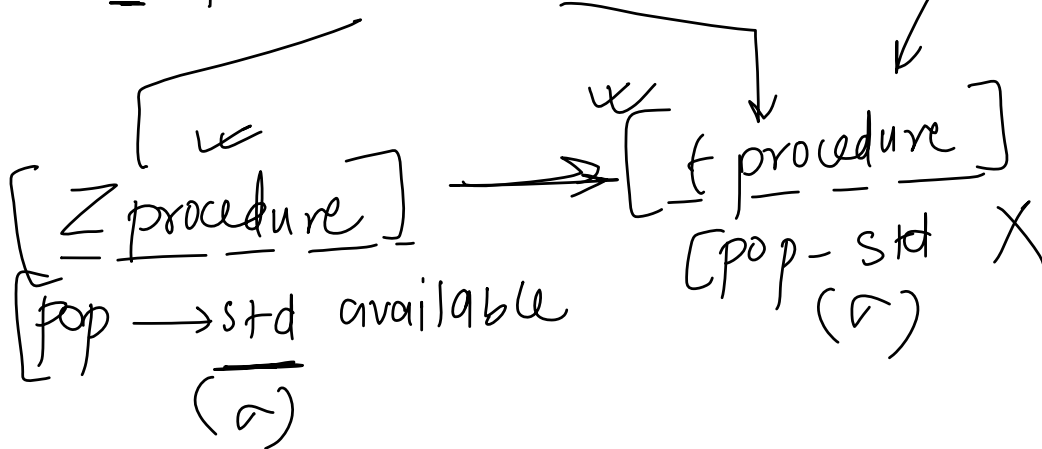
$95\%$ confident

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

$\rightarrow 95\%$

$25 \pm 4$    $[21, 29]$

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

Ways to calculate CI:

$[Z \text{ procedure}]$

$[pop \rightarrow std \text{ available}]$  $(\sigma)$

$[t \text{ procedure}]$

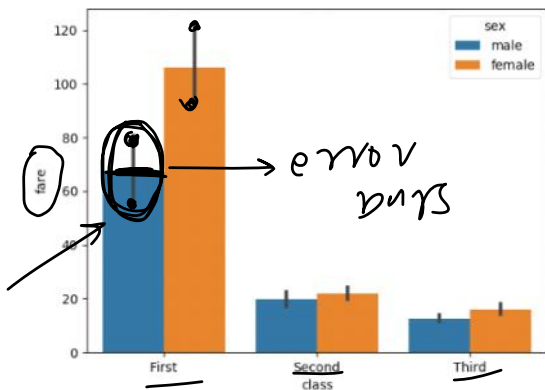$[pop - std \quad X$  $(\sigma)$

Confidence Interval is created for (Parameters) and not statistics. Statistics help us get the confidence interval for a parameter.

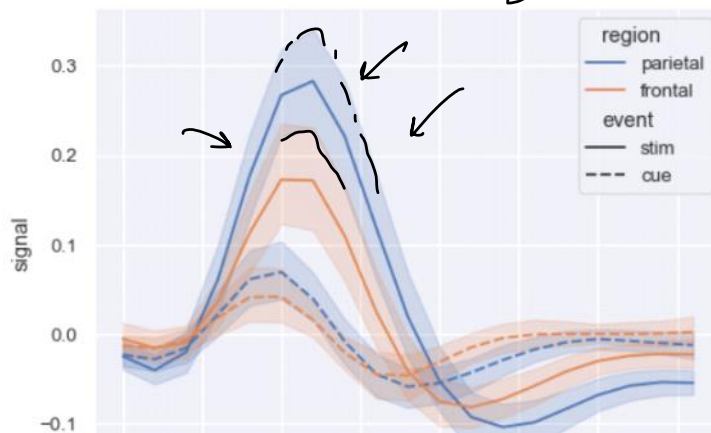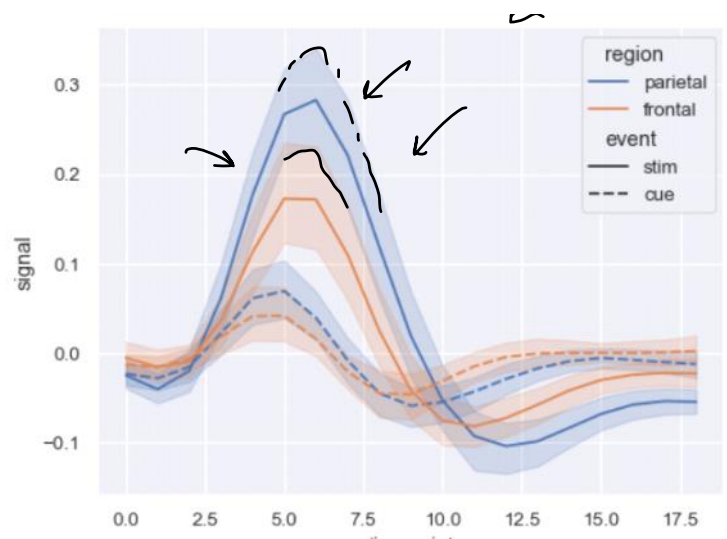[Examples of CT usage]    seaborn    $\rightarrow$ bar plot



error bars

(CI)

# Confidence Interval (Sigma Known)  Z procedure

30 March 2023    07:13

$(\sigma)$ pop std available → pop → sample random

YT → pop → std $(\sigma)$
distribution → normal
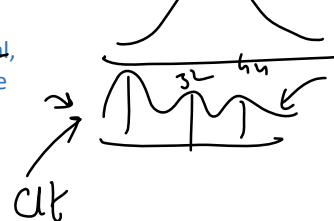
77k subs → age
→ normal dist

## Assumptions

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

2. Known population standard deviation: The population standard deviation (σ) must be known or accurately estimated. In practice, the population standard deviation is often unknown, and the sample standard deviation (s) is used as an estimate. However, if the sample size is large enough, the sample standard deviation can provide a reasonably accurate approximation.

3. Normal distribution or large sample size: The Z-procedure assumes that the underlying population is normally distributed. However, if the population distribution is not normal, the Central Limit Theorem can be applied when the sample size is large (usually, sample size n ≥ 30 is considered large enough). According to the Central Limit Theorem, the sampling distribution of the sample mean will approach a normal distribution as the sample size increases, regardless of the shape of the population distribution.

Sample size $\leq 30$  (20) → z procedure

CLT

A (1 - alpha)*100% Confidence Interval for mu:

$\sigma = 15$

YT → campusx → 77k → 28 ± 14 →

→ [16, 42] ← confidence interval

Confidence level → 95%

formula
CI using
Z procedure

$$CI = \overline{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

1) Intuition
2) $Z_{\alpha/2}$

$Z \to ?$

$(1-\alpha) \to$ confidence level

$(1-\alpha) \to 95\%$

$\sigma \to$ std pop

$n \to$ sample size → 100

Intuition
point estimate $(\overline{X}) \to$ CLT

point estimate $(\bar{x}) \to$ CLT

10 live class $\to$ (50) $\to$ avg class

$[\overline{x_1}, \overline{x_2} \cdots - \overline{x_{10}}] \to$ sampling dist of sample mean

$(z) \to N(0,1)$

normally dist



$\mu$

$\sigma/\sqrt{n}$

$(Z) = \dfrac{\bar{x} - \mu}{\sigma/\sqrt{n}}$

$95\% = (1-\alpha)$ sure

$\alpha = 5$

$CI = \bar{x} \pm (Z_{\alpha/2}) \dfrac{\sigma}{\sqrt{n}}$

$\boxed{95\%}$ $(z)$

$\boxed{1-\alpha}$ prob $\to 95\%$

$1-\alpha$

$\alpha/2$

Z live

$-Z_{\alpha/2}$ $Z_{\alpha/2}$

$-1$ $\quad 0$

$-Z_{2.5} \longleftrightarrow Z_{2.5}$

$\boxed{95\%}$

$\boxed{95}$

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

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

$(\mu) \to CI$

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

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

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

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

fixed → 95%

$\overline{X} \to$ Sample

$1-\alpha \quad 0.95 \quad \alpha = 0.\frac{50}{2}$

confidence $(1-\alpha) \to$ 95%

0.250

$$CI \qquad \underline{\mu} = \overline{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

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

$$\begin{array}{c} 950 \\ 025 \\ \hline 975 \end{array}$$

0.025    0.95    0.025

1.96

−1.96

range → confidence interval

$$\mu = \overline{X} \pm \boxed{1.96} \frac{\sigma}{\sqrt{n}} \longrightarrow \text{CI with a confidence in } -95\%$$

$Z_{\alpha/2} \to 1.96$    50%    75%    99%    $Z_{\alpha/2}$
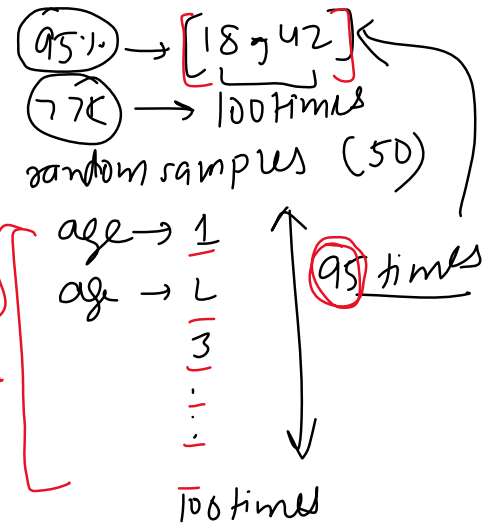
0.87

0.125    0.75    0.125

$Z \to 1.13$
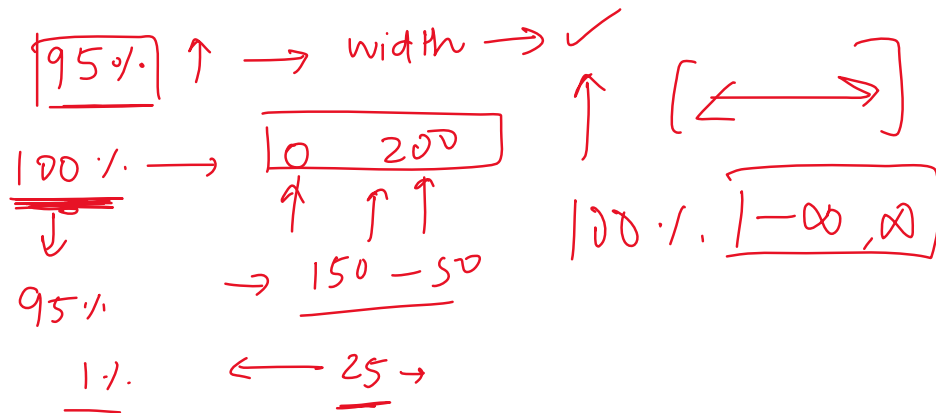
# Interpreting Confidence Interval

30 March 2023    08:33

*pop*
*fixed*    $45 \to [14-42] \leftarrow$

A confidence interval is a range of values within which a population parameter, such as the population 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 level**: The confidence level (commonly set at 90%, 95%, or 99%) represents the probability that the confidence interval will contain the true population 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 population and calculate the confidence interval for each, approximately 95 of those intervals would contain the true population parameter.

2. **Interval range**: The width of the confidence interval gives an indication of the precision of the estimate. A narrower confidence interval suggests a more precise estimate of the population 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 "X% confident that the true population 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 chose when constructing the interval.

What is the trade-off

*(handwritten annotations on right side:)*

$95\% \to [18, 42] \leftarrow$

$77\% \to 100 \text{ times}$

random samples (50)

age $\to 1$
age $\to 2$
$3$
$\vdots$

95 times

100 times

*(handwritten annotations at bottom:)*

$95\% \uparrow \to \text{width} \to \checkmark$

$100\% \to [0 \quad 200]$
$\downarrow$
$150 - 50$

$95\%$

$1\% \leftarrow 25 \to$

$\uparrow [\leftarrow \to]$

$100\% \quad [-\infty, \infty]$

# Factors Affecting Margin of Error

30 March 2023          07:15

Confidence ↑ → CI↑ → Margin of error↑

$$upper - lower$$ → marg.

1. Confidence Level (1-alpha) → $Z_{\alpha/2}$
2. Sample Size → more data → accurate prediction → $Z$ → CI↓
3. Population Standard Deviation

measure of variability ↑ → uncertainty↑ → CI↑

$$CI = point \ estimate \pm \boxed{margin \ of \ error}$$

$$= \boxed{\bar{x}} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

→ pop std ②
sample size ③

Sample mean

① critical value



Relationship Between Margin of Error and Critical Value (Z Procedure)

∞

95

# Confidence Interval (Sigma not known)

$\boxed{\sigma} \longrightarrow$ t-procedure

Using the t procedure

Assumptions

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

2. Sample standard deviation: *majboori* The population standard deviation (σ) is unknown, and the sample standard deviation (s) is used as an estimate. The t-distribution is specifically designed to account for the additional uncertainty introduced by using the sample standard deviation instead of the population standard deviation.

3. Approximately normal distribution: The t-procedure assumes that the underlying population is approximately normally distributed, or the sample size is large enough for the Central Limit Theorem to apply. If the population distribution is heavily skewed or has extreme outliers, the t-procedure may not be accurate, and non-parametric methods should be considered.

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

Sample $> 30 \rightarrow$ If it's not normal

normal $\rightarrow$ small sample size $\rightarrow$ CI

$t_{\alpha/2} \gg z_{\alpha/2}$

$t_{\alpha/2} \simeq z_{\alpha/2}$

z-table

$$CI = \overline{X} \pm z_{\alpha/2} \frac{\boxed{\sigma}^{X \longrightarrow S}}{\sqrt{n}}$$

pop mean

$\overline{X}$

complexity

$$CI = \overline{X} + \boxed{z_{\alpha/2}} \frac{S}{\sqrt{n}}$$

$$CI = \overline{X} + \boxed{t_{\alpha/2}} \frac{\boxed{S}}{\sqrt{n}} \quad \boxed{ds} > 30$$

$95\% \; z_{\alpha/2} \rightarrow 1.96 \quad t_{\alpha/2} = 2.0008$

$\boxed{z}$

$$\boxed{\frac{\overline{X} - \mu}{\sigma/\sqrt{n}}} \rightarrow \text{standard normal dist}$$

Student's t distribution

$n-1$

$\boxed{z'}$

$$\boxed{\frac{\overline{X} - \mu}{\boxed{S}/\sqrt{n}}}$$

param Uu

$\boxed{df} \rightarrow \boxed{n-1}$

theoretical dist

$\boxed{49}$  $\boxed{30}$

$\rightarrow$ degree of freedom

In student t distribution $\rightarrow$

$\rightarrow n \uparrow \rightarrow \boxed{t = N}$

Student's t dis

$df \rightarrow \infty \; (t \rightarrow N)$ ----

$\mu, \sigma$

Because of theoretical dist $\rightarrow$ uncertainty $\uparrow$

$\rightarrow$ thick tail

than CI

for right inference

# Normal and t-Distribution Comparison



Probability Density vs x

- Normal Distribution (blue)
- Student's t-Distribution (df=6) (red)

0.950

2.5        2.5
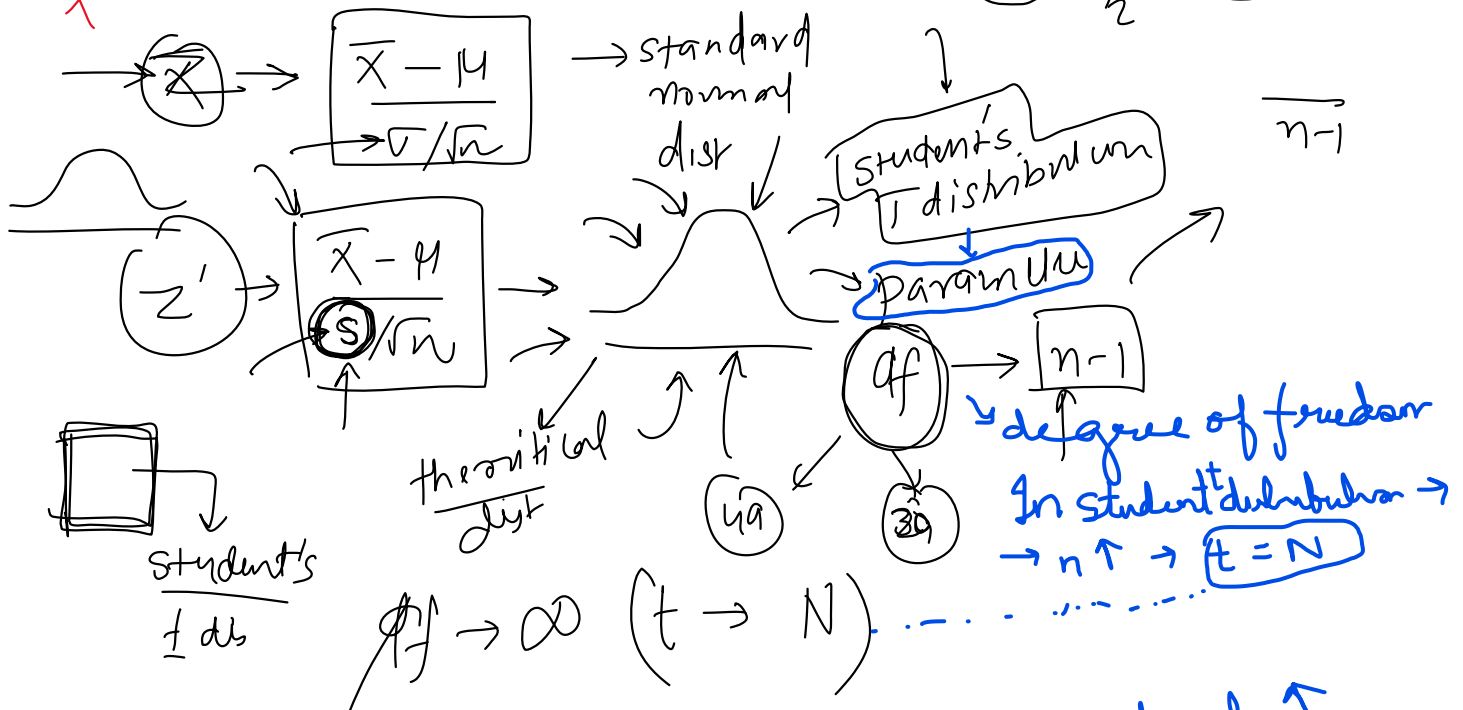
$$CI = \bar{X} \pm Z_{\alpha/2} \frac{S}{\sqrt{n}}$$

$\sigma$

$$= \bar{X} \pm t_{\alpha/2} \frac{S}{\sqrt{n}}$$

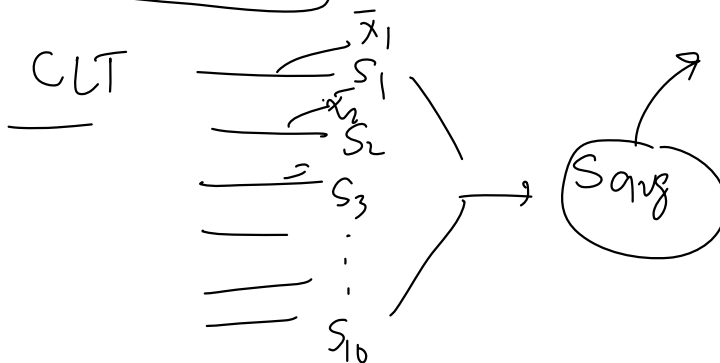$$CI = \bar{X} \pm t_{\alpha/2} \frac{S}{\sqrt{n}} \longleftrightarrow CI = \bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

CLT

$\bar{X}_1$
$S_1$
$\bar{X}_2$
$S_2$
$S_3$
$\vdots$
$S_{10}$

$S_{avg}$

# Student's T Distribution

30 March 2023    07:16

Student's t-distribution, or simply the t-distribution, is a probability distribution that arises when estimating the mean of a normally distributed population when the sample size is small and the population standard deviation is unknown. It was introduced by William Sealy Gosset, who published under the pseudonym "Student."

The t-distribution is similar to the normal distribution (also known as the Gaussian distribution or the bell curve) but has heavier tails. The shape of the t-distribution is determined by the degrees of freedom, which is closely related to the sample size (degrees of freedom = sample size - 1). As the degrees of freedom increase (i.e., as the sample size increases), the t-distribution approaches the normal distribution.

In hypothesis testing and confidence interval estimation, the t-distribution is used in place of the normal distribution when the sample size is small (usually less than 30) and the population standard deviation is unknown. The t-distribution accounts for the additional uncertainty that arises from estimating the population standard deviation using the sample standard deviation.

To use the t-distribution in practice, you look up critical t-values from a t-distribution table, which provides values corresponding to specific degrees of freedom and confidence levels (e.g., 95% confidence). These critical t-values are then used to calculate confidence intervals or perform hypothesis tests.

# Titanic Case Study

$$POP \rightarrow 1360$$

$$\mu \rightarrow X$$

$$\sigma \rightarrow X$$

$$CLT \rightarrow 10 \text{ times} \rightarrow \underline{30} \text{ size}$$

95% confidence level

inferenc