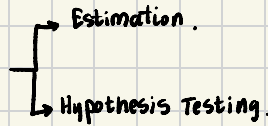


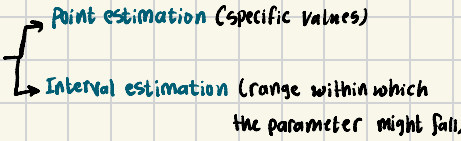
Chapter 6 -

Estimation

Inferential statistics



Estimation



population & sample

Reference / Target / Study population

This is where we obtain our samples from.

↳ The Group we want to study.

populations



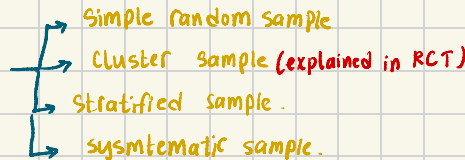
not every "large" population is infinite

Sample → Some members of a population.

a Random sample :- the sample is random if each member is chosen independently and has a known (non-zero) probability of being chosen

A Simple Random Sample → is a Random sample where each member has an equal chance of being selected.

Types of probability Random samples

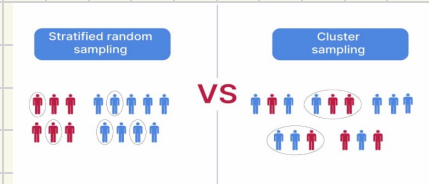


↳ cluster sampling →

when geographic, spatial characteristics naturally divide the target population into random groups (clusters) like schools, neighborhoods

↳ Stratification :-

method of dividing a target population into subgroups (strata) based on characteristics thought to be important (age, gender, clinical condition)



Notation:-

The probability distribution of a Sample Statistic is called the sampling distribution
↳ like: \bar{x} , \hat{p} , s^2

① estimation of \bar{x}

$\hat{\theta} \rightarrow$ sample statistic

$\theta \rightarrow$ population parameter

If: $E(\hat{\theta}) = \theta$

then $\hat{\theta}$ is a good estimator of θ

\bar{x} is a good estimator of μ because \rightarrow

$\triangleright E(\bar{x}) = \mu \rightarrow$ population mean

$\triangleright \bar{x}$ (sample mean) is called **Minimum variance unbiased estimator (more accurate)**

Notation \rightarrow

even when using a good estimator like \bar{x} , the precision of the estimation can still be affected by

\triangleright Standard Error of Sample mean \rightarrow

The variance (σ^2) of the sample mean

$\text{var}(\bar{x}) = \frac{\sigma^2}{n} \Rightarrow$ the standard deviation of \bar{x} is called the **standard error** $\left[\frac{\sigma}{\sqrt{n}}\right]$

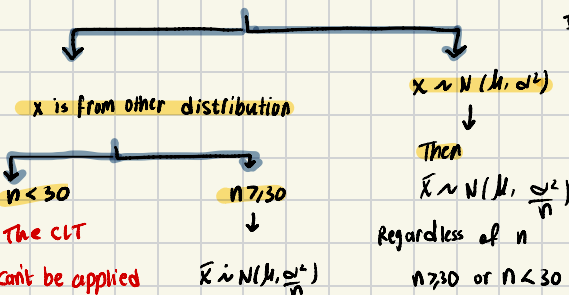
the standard error can be approximated by: -

$\frac{s}{\sqrt{n}}$
 \nearrow sample standard deviation

\hookrightarrow when σ is unknown.

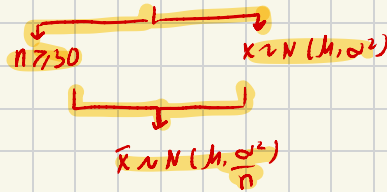
\triangleright The Central Limit theorem

if



• Standardization of \bar{x}

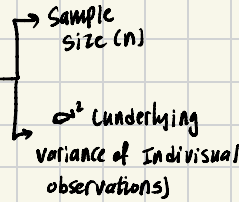
whether



Z score of \bar{x} can be found using :-

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$\sim N(0, 1)$
 Standard normal distribution.



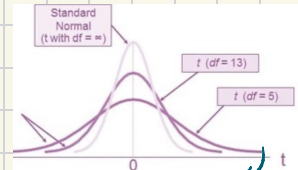
t-distribution

- bell shaped (symmetric)
- has mean = 0 \rightarrow greater than 1
- variance = $\frac{df}{df-2}$, $df \rightarrow$ degrees of freedom $df = n-1$
 \uparrow sample size

- The t distribution has "thicker tails" [higher variance] than the standard normal distribution

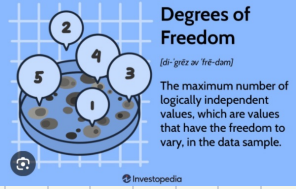
But as the sample size (n) increases the variance decreases such that $df \approx \infty \Rightarrow$ The t

distribution is the same as the standard normal.



$n <$ then variance $>$

degrees of freedom

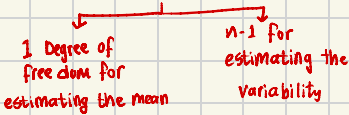


example on Degrees of freedom

Consider a data sample consisting of five positive integers. The values of the five integers must have an average of six. If four items within the data set are (3, 8, 5, and 4), the fifth number must be 10. Because the first four numbers can be chosen at random, the degree of freedom is four.

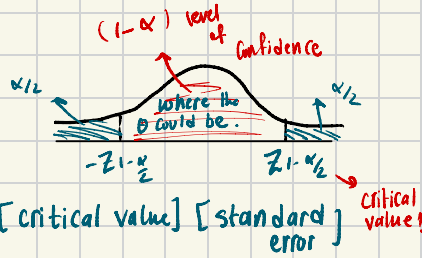
Degrees of freedom = $n - 1$

In the case of the t-distribution, the degrees of freedom are $N - 1$ as one degree of freedom is reserved for estimating the mean, and $N - 1$ degrees remain for estimating the variability.



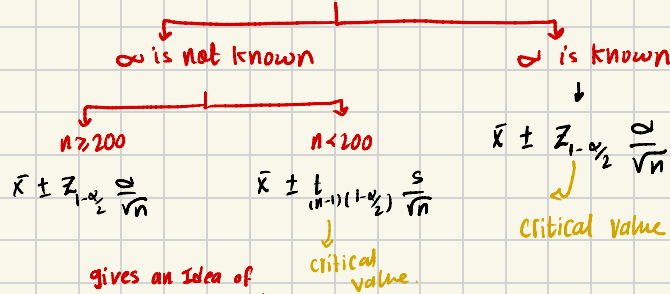
interval estimation of mean

Confidence Interval
- level of Confidence $(1 - \alpha)$



CI = point estimator \pm [critical value] [standard error]

to construct a $(1 - \alpha) 100\%$ CI for (μ)



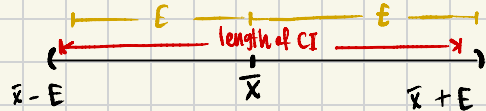
gives an idea of the precision of the point estimate

meaning of $(1 - \alpha) 100\%$ CI for μ .

Over the collection of all 95% CIs that could be constructed from repeated random samples of size n , 95% will contain the parameter μ .

length of CI

length of Confidence interval = upper bound - lower bound



CI = point estimator \pm (critical value) (standard error)

length of CI = $2E$

example $\rightarrow L = 2 \times Z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}$

length

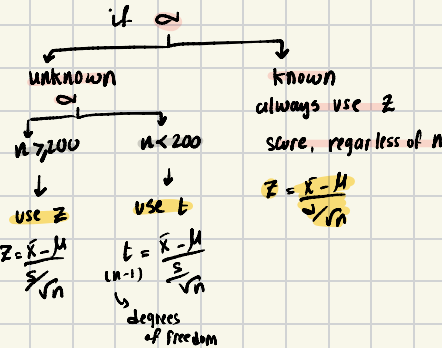
factors that affect the length of CI

- 1 level of Confidence $(1 - \alpha)$ directly
- 2 $\alpha \Rightarrow$ Inversly
- 3 sample size (n) Inversly
- 4 standard error directly

$\frac{\sigma}{\sqrt{n}}$ \rightarrow in that case S or length of CI directly

length $\propto (1 - \alpha) \propto \frac{1}{\alpha} \propto S \propto \frac{1}{n}$

when to use Z or t?



notation

when $n \geq 30$

the t distribution variance is nearly (1)
so there won't be a huge difference between the Z of t distribution!

So Don't worry about that please 😊

► Estimation for the binomial distribution

$$X \sim \text{Bin}(n, p)$$

◀ population proportion

► \hat{p} (sample proportion)

$$\hat{p} = \frac{X}{n} \begin{array}{l} \rightarrow \text{number of successes} \\ \rightarrow \text{sample size} \end{array}$$

► The sample proportion \hat{p} is a good point estimator of p

- $E(\hat{p}) = E\left(\frac{X}{n}\right) = \frac{1}{n} E(X) = \frac{np}{n} = p$
- $\text{Var}(\hat{p}) = \text{Var}\left(\frac{X}{n}\right) = \frac{1}{n^2} \text{Var}(X) = \frac{1}{n^2} (npq) = \frac{pq}{n} \approx \frac{\hat{p}\hat{q}}{n}$ if p is unknown
- standard error of $\hat{p} = \frac{pq}{n} \approx \frac{\hat{p}\hat{q}}{n}$ if p is unknown.

from the central limit theorem

$\hat{p} = \bar{X}$ is normally distributed

$$\hat{p} \sim N\left(p, \frac{pq}{n}\right), \text{ as long as } n\hat{p}\hat{q} \geq 5$$

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}}$$

confidence interval of binomial proportion

$$\hat{p} \pm \sqrt{\frac{\hat{p}\hat{q}}{n}} Z_{1-\alpha/2} \quad n\hat{p}\hat{q} \geq 5$$

↳ maximum standard error ↑,
↑ length is when $\hat{p} = \hat{q} = \frac{1}{2}$