

Unit - I

sampling distribution

Inferential statistics

It is a statistical technique to draw information about population using information based on sample only.

Two branches

1. Estimation of parameter

2. Hypothesis testing

- parametric test

- Non-parametric test

Key concept

1. Sampling distribution

2. Probability distribution

3. Sampling

What is sampling distribution?

→ The concept of sampling distribution is very important in inferential because almost all inferential statistics are based on sampling distribution.

It is the probability distribution of all possible values that a statistics can assume (say sample mean, sample proportion, etc) computed from sample of same size, drawn randomly and repeatedly from same population.

1. Sampling distribution of mean
2. Sampling distribution of proportion.
3. Sampling distribution of sample standard deviation
4. Sampling distribution of sample correlation coeff.

Three measures describing a sampling distribution

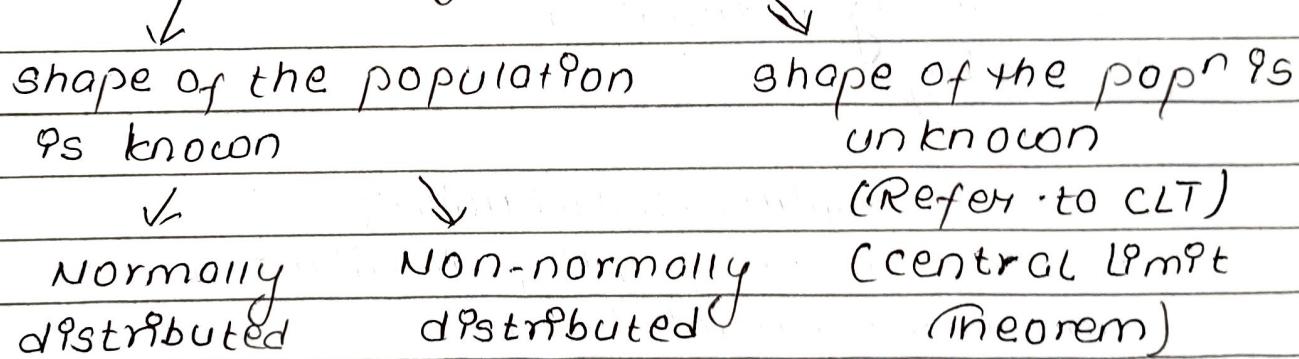
1. Fractional form (probability distribution of stat.)
2. Expected value (Average value of stat.)
3. Standard Error (standard deviation of stat.)

sampling distribution of sample Mean (\bar{X})

1) Functional Form

(Probability distribution of \bar{X})

Three sampling situation



what sample size is large enough for normal approximation of sampling distribution of mean (\bar{X})?

- 1. For the most population distribution, regardless of the shape, the sampling distribution of the mean (\bar{X}) is approximately normally distributed if sample of at least 30 observations are selected.

$n \geq 30$ (large sample)

$n < 30$ (small sample)

2. If the population distribution is fairly symmetrical, the sampling distribution of the mean is approximately normally distributed if sample of at least 15 obs. are selected.

3. If the population is normally distributed, the sampling distribution of mean is normally distributed regardless of the sample size.

2) Expected value

For any sample of 'n' observations, the expected value of the sample mean is equal to population mean.

$$E(\bar{X}) = \text{Mean of the sample mean}$$

$$= \mu_{\bar{X}}$$

$$= \mu$$

This is true for both SRSWR and SRWSWOR.

$$E(\bar{X}) = \mu$$

→ The sample mean (\bar{X}) is an unbiased estimation of population mean (μ).

→ The sample mean will be neither too large nor too small and it is close to μ as long as n is large.

3) Standard error

The standard deviation of the sample mean is

called its standard error. It is the measure of variability among the sample means. Higher the standard error, greater is the variability among the sample mean.

$$S.E.(\bar{X})_{SRSWOR} = \frac{\sigma}{\sqrt{n}} \quad [\sigma = \text{population s.d.}] \quad [n = \text{sample size}]$$

6 knowledge

1. From previous related large scale study.
2. Meta Analysis.

$$3. \sigma \approx \frac{R}{6}$$

$$\text{Estimated } S.E.(\bar{X}) = \frac{s}{\sqrt{n}} \quad [s = \text{sample standard dev.}]$$

Factors affecting standard error.

1. variability of characteristics in the population (Direct)
2. sample size (Indirect)

$$S.E.(\bar{X})_{SRSWOR} = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} \quad [\text{is called FPC}] \quad [\text{finite pop'n correction}]$$

$$\text{Estimated } S.E.(\bar{X})_{SRSWOR} = \frac{s}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

First Formula

1. Infinite popn
2. SRSWOR
3. $f < 5\%$ ($f = n/N$ = sampling fraction)

Second Formula

1. Finite popn
2. SRSWOR
3. $f \geq 5\%$

known \downarrow unknown
 $f < 5\% \leftarrow \rightarrow f \geq 5\%$.

Central ImpPt Theorem (CLT)

For the case where sampling is from a non-normally distributed population, we refer to an important mathematical theorem called CLT.

Statement

As the sample size (i.e. no. of observations in each sample) increases, the sampling distribution of the mean will approach normality. This is true regardless of the shape of the distribution of the individual values in the population.

SRSWOR \rightarrow Nⁿ
SRSWOR \rightarrow N_n



Numerical Example

- # The following data represent the no. of days of absent per year in a population of employees of a small IT company.

8, 3, 1, 11, 4, 7

- (a) consider all possible samples of size two, i.e. $n=2$, which can be drawn without replacement.
- (b) Draw histogram of distribution of population values and sampling distribution of means. Comment on the shape of distribution.
- (c) Find the mean of the population and mean of the samples and verify $E(\bar{X}) = \mu$.
- (d) Find the population standard deviation and verify that $S.E(\bar{X}) = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$.

\rightarrow

Solution:

(a) X = No. of days of absent in a year

N = population size = 6

n = sampling size = 2

Number of possible samples = $N_{n,n}$ (SRSWOR)
 $= 6 C_2$
 $= 15$

Possible samples

Sample No.	Sample	Sample Mean
1	8, 3	5.5
2	8, 1	4.5
3	8, 11	9.5
4	8, 4	6
5	8, 7	7.5

6	3, 1	2
7	3, 11	7
8	3, 4	3.5
9	3, 7	5
10	1, 11	6
11	1, 4	2.5
12	1, 7	4
13	1, 14	7.5
14	1, 7	9
15	4, 7	5.5

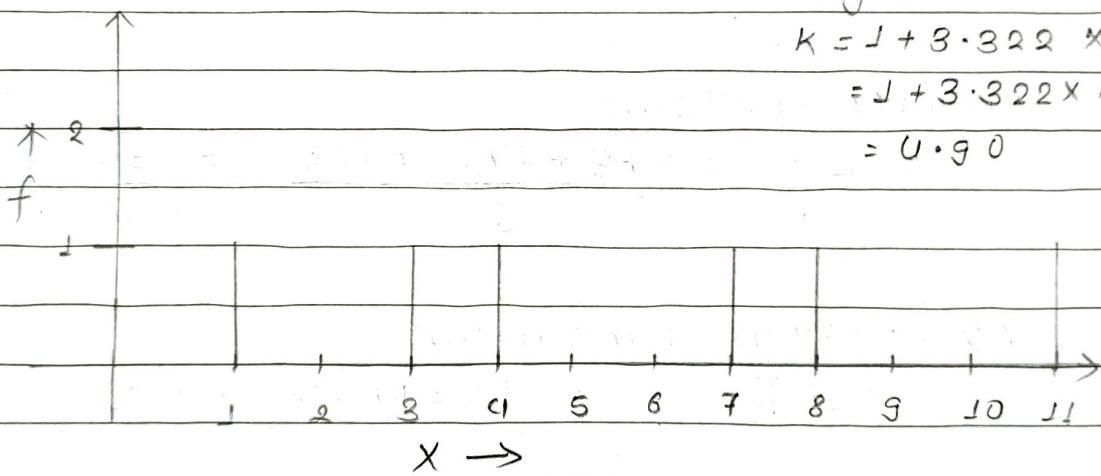
(b) shape of the population data

$$\text{Range} = 9.5 - 2 = 7.5$$

$$K = 1 + 3.322 \times \log n$$

$$= 1 + 3.322 \times \log 15$$

$$= 6.90$$



shape of the popn is rather uniform

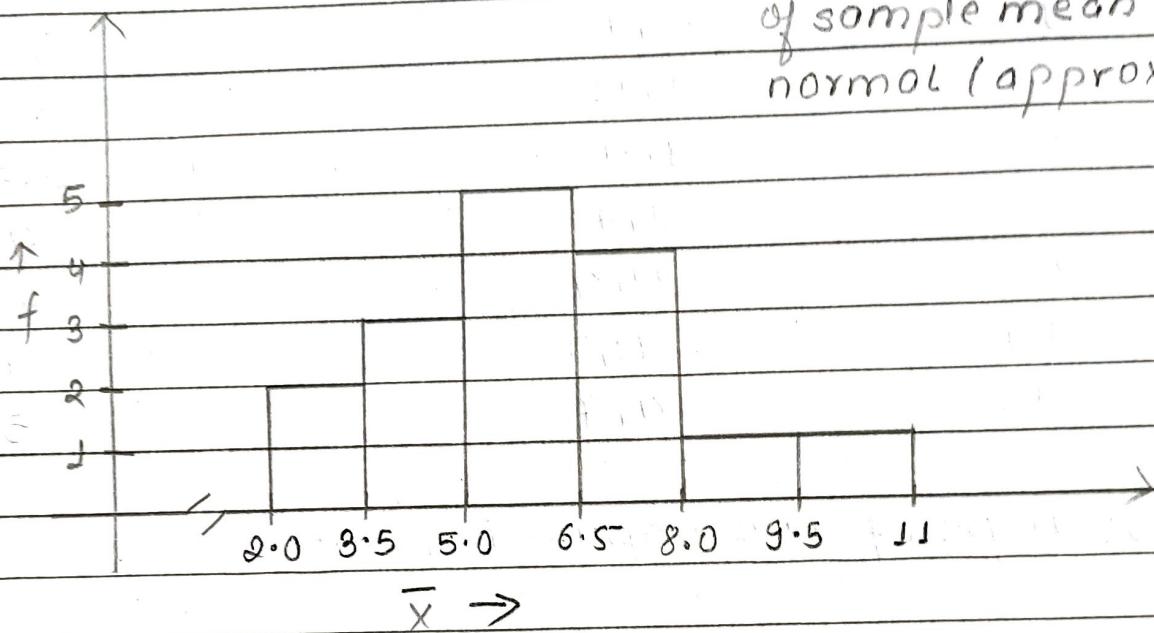
(No. of samples)

Sample Mean Frequency

2.0 - 3.5	2	→
3.5 - 5.0	3	
5.0 - 6.5	5	
6.5 - 8.0	3	
8.0 - 9.5	1	
9.5 - 11	1	

Shape of the sampling distribution of mean

shape of the distribution
of sample mean is rather
normal (approximately)



② Population Mean

$$\mu = \frac{\sum x}{N} = \frac{8+3+1+11+4+7}{6} = \frac{34}{6} = 5.67$$

$E(\bar{x})$ = Mean of the mean

$$= \frac{\sum \bar{x}}{K} \quad [K = \text{No. of samples} = N c_n = 15]$$

$$\begin{aligned} &= 5.5 + 4.5 + 9.5 + 6 + 7.5 + 2 + 7 + \dots + 5.5 \\ &\qquad\qquad\qquad 15 \\ &= 5.67 \end{aligned}$$

$$\therefore E(\bar{x}) = \mu$$

verified.

$$\textcircled{d} \quad S.E(\bar{x}) = \frac{6}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

population SD:

$$\begin{aligned} \sigma &= \sqrt{\frac{\sum (x - u)^2}{N}} = \sqrt{\frac{\sum x^2 - u^2}{N}} \\ &= \sqrt{\frac{8^2 + 3^2 + 1^2 + 11^2 + 4^2 + 7^2 - (5.67)^2}{6}} \\ &= \sqrt{\frac{180 - 32.1489}{3}} \\ &= 3.3443 \end{aligned}$$

$$\begin{aligned} \text{RHS} &= \frac{6}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} \\ &= \frac{3.3443}{\sqrt{2}} \sqrt{\frac{6-2}{6-1}} \\ &= 2.1151 \end{aligned}$$

$$\begin{aligned} \text{LHS} &= S.E(\bar{x}) \\ &= \sqrt{\frac{\sum (\bar{x} - \bar{\bar{x}})^2}{K}} \end{aligned}$$

$$\begin{aligned} \text{LHS} &= S.E(\bar{x}) \\ &= \sqrt{\frac{\sum (\bar{x} - \bar{\bar{x}})^2}{K}} \quad \left| \begin{array}{l} \bar{x} = \text{Mean of the mean} \\ K = \text{No. of samples} = 15 \end{array} \right. \\ &= \sqrt{\frac{(5.5 - 5.67)^2 + (4.5 - 5.67)^2 + \dots + (5.5 - 5.67)^2}{15}} \\ &= \sqrt{\frac{0.0289 + 1.3689 + 14.6689 + 0.1089 + 3.3489 + 13.4689 + 1.7689 + 4.7089 + 0.4489 + 0.1089 + 10.0489 + 2.7889 + 3.3489 + 11.0889 + 0.0289}{15}} \\ &= \sqrt{\frac{67.3335}{15}} = 2.1187 \end{aligned}$$

verified

Sampling distribution of sample population (\hat{p})

Parameter $\rightarrow P$ or π (population proportion of success)
 Statistic $\rightarrow \hat{p}$ (sample proportion of success)

Consider sampling from a population having two classes or categories:

- One class possessing a particular attribute (success)
- other class not possessing a particular attribute (failure)

The presence of attribute in the sample unit is known as success, and its absence as failure.

Terminology

P or π = proportion of success in the population (parameter)

$$= \frac{X}{N} = \frac{\text{No. of success in population}}{\text{population size}}$$

$1 - P$ or $1 - \pi$ or Q = proportion of failure in the population (parameter)

$$= \frac{N-X}{N} = \frac{\text{No. of failures in the population}}{\text{population size}}$$

\hat{p} = proportion of success in the sample (statistics)

$$= \frac{X}{n} = \frac{\text{No. of successes in the sample}}{\text{sample size}}$$

for $1-p$ = proportion of failure in the sample

$$= \frac{n-x}{n}$$

= No. of failures in the sample
~ Sample size

1. Functional form

(Probability distribution of p)

When the sample size is large, the sampling distribution of sample proportion (p) is approximately normally distributed by virtue of central limit theorem.

How large?

- Both np and $n(1-p)$ must be greater than 5

2. Expected value of p

The sample proportion of success ' p ' is an unbiased estimation of population proportion of success ' p '.

$$\text{P.E. } E(p) = p$$

⇒ Mean of sample proportion of success.

= population proportion of success.

3. Standard error of p

The standard deviation of sample proportion of success ' p ', is called its standard error. It measures the variability of ' p ' in the repeated sampling.

$$1) S.E.(p) = \sigma_p = \sqrt{\frac{pxq}{n}} \quad [Infinite popn or SRSWR or f < 5\%]$$

$$2) S.E.(p) = \sigma_p = \sqrt{\frac{pxq}{n}} \sqrt{\frac{n-n}{N-1}} \quad [Finite popn or SRSWOR or f > 5\%]$$

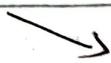
Population



Finite correction multipl.
user (FCM)

Infinite
(FCM X)

Finite



$f \geq 5\%$.

(FCM V)

$f < 5\%$.

(FCM X)

$$\text{Est. S.E.}(\hat{p}) = \sqrt{\frac{p \times q}{n}} \quad \text{case I}$$

$$\text{Est. S.E.}(\hat{p}) = \sqrt{\frac{p \times q}{n} \cdot \frac{n - n}{N - 1}} \quad \text{case II}$$

Internal estimation of population mean (μ)

when population standard deviation (σ) is known.

→ when population S.D. (σ) is known & the sampling distribution of sample mean (\bar{x}) follows normal probability distribution with mean (μ) & standard deviation of $\frac{\sigma}{\sqrt{n}}$

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$

The standardized value of \bar{x} is,

$$z = \frac{\bar{x} - E(\bar{x})}{SD(\bar{x})} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$P(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha \quad (1 - \alpha) \times 100\% \text{ C.I.}$$

The $(1 - \alpha) \times 100\%$ confidence interval for population mean is given by,

$$P\left\{-Z_{\alpha/2} \leq \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \leq +Z_{\alpha/2}\right\} = 1 - \alpha$$

$$= P\left\{-Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \bar{x} - \mu \leq +Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right\} = 1 - \alpha$$

$$= P\left\{-\bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq -\bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right\} = 1 - \alpha$$

$$= P\left\{\bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right\} = 1 - \alpha$$

This is called $(1 - \alpha) \times 100\%$ confidence interval for μ .

$$\hat{\mu}_L = \bar{x} - Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$$\hat{\mu}_U = \bar{x} + Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Then, $(1-\alpha) \times 100\%$ confidence interval for μ is given by,

$$\bar{x} + \frac{Z_{\frac{\alpha}{2}} \cdot \sigma}{\sqrt{n}}$$

Margin of error.

where,

\bar{x} = Sample mean

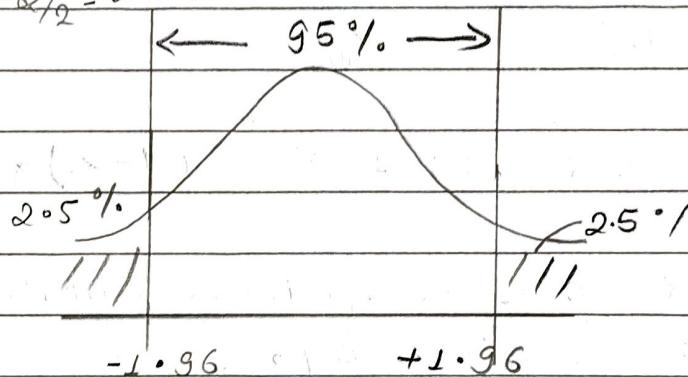
$Z_{\frac{\alpha}{2}}$ = Critical value of z for particular level of significance α

σ = Standard error of \bar{x}

$$\sqrt{n} \quad \alpha = 5\%$$

$$\alpha/2 = 2.5\%$$

$$1 - \alpha = 95\%$$



Notes

- 1. If population is finite, and $f = n/N > 5\%$, we have to use finite correction multiplier.

$$\bar{x} \pm Z_c \cdot \frac{\sigma}{\sqrt{n}} \cdot \sqrt{\frac{N-n}{N-1}}$$

when population standard deviation (σ) is unknown

→ when population standard deviation (σ) is unknown, the sampling distribution of \bar{x} follows student's t-distribution with $n-1$ degrees of freedom.

The $(1-\alpha) \times 100\%$, confidence interval for population mean (μ) is given by,

$$\bar{x} \pm t_{\frac{\alpha}{2}}(n-1) \cdot \frac{s}{\sqrt{n}}$$

where,

\bar{x} = Sample mean

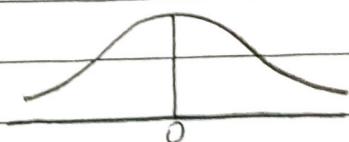
s = Sample S.D.

n = Sample size

$t_{\alpha/2}(n-1)$ = Critical value of t for α level of significance and $n-1$ degrees of freedom.

Characteristics of t-distribution

1. Shape : t-distribution is symmetrical distribution with mean = 0 at its center.



2. The t-distribution is exact prob-distribution. There is a particular t-curve for particular sample size.

3. As the sample size increases the t-distribution approaches z-distribution.

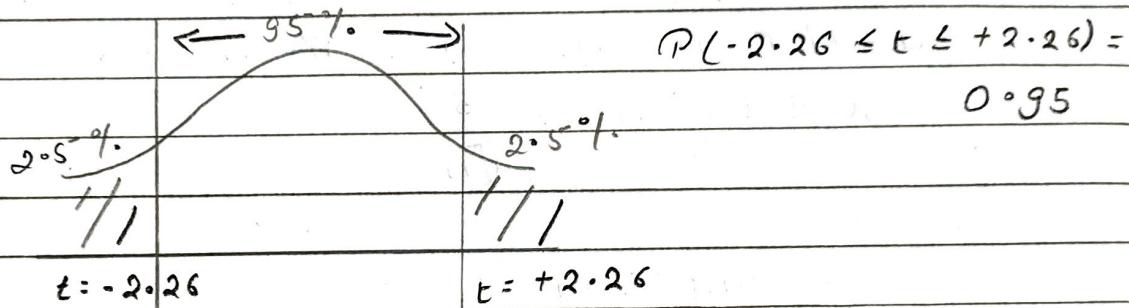
4. To find the t-value we need information on:

(i) α -value

(ii) degree of freedom

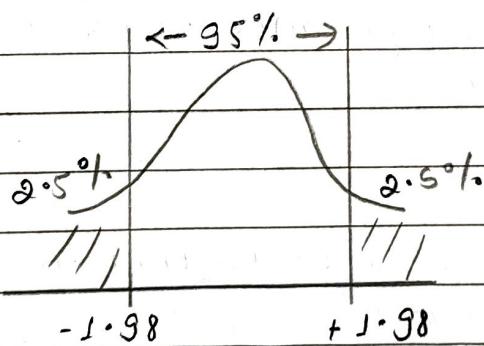
$$\textcircled{i} \quad \alpha = 5\% , n = 10 , \text{degree of freedom} = 10 - 1 = 9$$

$$t_c = \pm 2.26$$



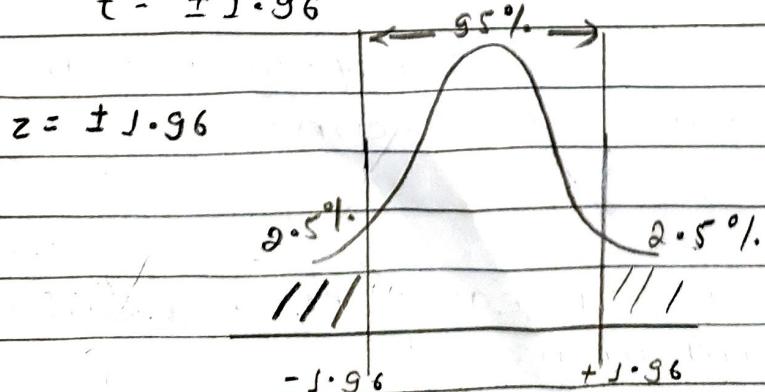
$$\textcircled{ii} \quad \alpha = 5\% , n = 100 , \text{d.f.} = 100 - 1 = 99$$

$$t_c = \pm 1.98$$



$$\textcircled{iii} \quad \alpha = 5\% , n = 1000 , \text{d.f.} = 1000 - 1 = 999$$

$$t_c = \pm 1.96$$



Note:

1. If population is finite and $f > 5\%$, use finite correction multiplier.

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

2. If we take large sample, we can replace t-value with z-value

Numerical:

Interval Estimation (confidence interval) for population mean (μ) when population standard deviation (σ) is known.

- # An electrical firm manufactures 10 light bulbs that have a length of life that is approximately normally distributed with the standard deviation of 40 hours. If a sample of 30 light bulbs has an average life of 780 hrs, find the 95% confidence interval for the population mean of all light bulbs produced by this firm.

→ Solution:

X = length of life of light bulbs (hrs)
 X has normal probability distribution.

Given,

σ = Population S.D. = 40 hours.

n = Sample size, i.e. no. of test bulbs in the sample = 30

\bar{X} = Sample mean life of light bulbs = 780 hours.

α = Significance level / probability = 5%.

$1 - \alpha$ = confidence level = 95%.

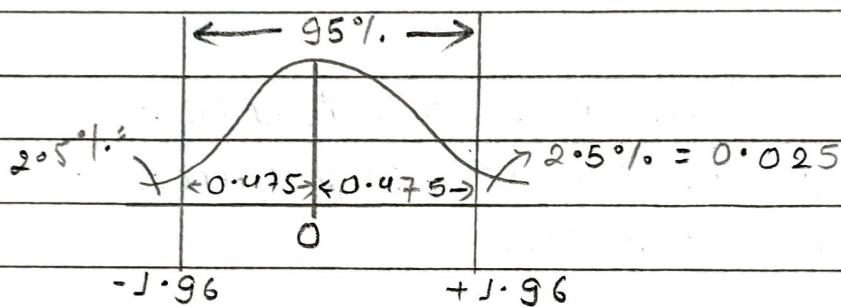
$$\alpha = ?$$

Now,

The $(1 - \alpha) \times 100\%$ confidence level interval for μ is given by,

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

Critical z ,



$$\text{critical } z = z_c$$

$$= Z_{0.025}$$

$$= 1.96$$

Now, 95% confidence interval for population (μ) is given by,

$$780 \pm 1.96 \cdot \frac{40}{\sqrt{30}}$$

$$= 780 \pm 14.3138$$

$$= 794.3138, 765.6861$$

Conclusion:

We are 95% sure that the true average life of the light bulb produced by the firm is between 765.7 hrs to 794.3 hrs.

confidence interval of μ using t-distribution
(σ unknown)

The following measurements were recorded for the drying time, in hours of a certain brand of latex point.

3.4	2.5	4.8	2.9	3.6
2.8	3.3	5.6	3.7	2.8
4.9	4.0	5.2	3.0	4.8

Assuming that the measurements represent a random sample from a normal population, find the 99% confidence interval for the mean drying time.

→ Solution:

X = Drying time of latex point (hrs)
 X has normal probability distribution.

Given,

n = Sample size

= No. of test points in the sample = 15

σ = popn of SD = ?

The $(1-\alpha) \times 100\%$. confidence interval for μ is given by,

$$\bar{x} \pm t_{\alpha/2}(n-1) \cdot \frac{s}{\sqrt{n}} \quad \begin{array}{l} \text{choice of distn in t} \\ \text{because } \sigma \text{ is unknown} \end{array}$$

Now,

$$\begin{aligned}\sum X &= 3.4 + 2.5 + \dots + 4.8 \\ &= 56.8\end{aligned}$$

$$\begin{aligned}\sum X^2 &= 3.4^2 + 2.5^2 + \dots + 4.8^2 \\ &= 228.28\end{aligned}$$

$$\bar{x} = \frac{\sum x}{n} = \frac{56.8}{15} = 3.7867$$

$$s = \sqrt{\frac{1}{n-1} \left\{ \sum x^2 - n \cdot \bar{x}^2 \right\}}$$

$$= \sqrt{\frac{1}{15-1} \left\{ 228.28 - 15 \times (3.7867)^2 \right\}}$$

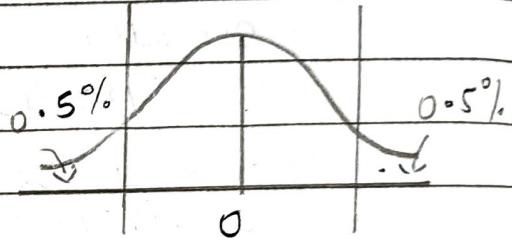
$$= 0.9708$$

Critical

α = significance level = 1%

$1-\alpha$ = confidence level = 99%

d.f = degrees of freedom
 $= n - 1 = 14$



From table,

$$t_c = 2.977$$

∴ The 99% confidence interval for μ given by,

$$56.8 \pm 2.977 \times \frac{0.9708}{\sqrt{15}}$$

$$= 56.8 \pm 0.7462$$

$$= 57.5462, 56.0537$$

Conclusion:

We are 99% sure that the true average drying time is between 56.0537 and 57.5462 hrs.

Interval estimation of population proportion of success (π)

→ If we took large sample ' n ', the sampling distribution of sample proportion of success ' p ' has normal distribution, with mean of π (or P) and standard deviation of $\sqrt{\frac{pxq}{n}}$.

The $(1-\alpha) \times 100\%$ confidence interval for π is given by $p \pm z_c \sqrt{\frac{pxq}{n}}$

where, p = sample proportion of success.

q = sample proportion of failure.

n = sample size

z_c = critical Z .

A survey of 500 people shopping mall, selected at random, showed that 350 of them used credit cards for their purchases and 150 used cash. Construct a 95% confidence interval estimate of proportion of all persons at the mall, who used credit card for shopping.

→ Solution:

x = Mode of payment (credit card, cash)

The population is has two categories:

- success (using credit card for payment)

- failure (using cash)

n = sample size

= No. of people in the survey = 500.

p = proportion of success in sample

$$= \frac{X}{n}$$

[X = no. of successes in sample]

$$= \frac{350}{500} = 0.70$$

q = proportion of failure in sample

$$= \frac{n-X}{n} = \frac{150}{500} = 0.3 \quad | \quad q = 1-p \\ = 1 - 0.7 = 0.3$$

The $(1-\alpha) \times 100\%$ confidence interval for true proportion of success (π) is given by,

$$P \pm Z_c \sqrt{\frac{pq}{n}}$$

$$Z_c = \text{critical } Z \text{ for } \alpha = 5\% \\ = 1.96$$

$$\alpha = \text{level of significance} \\ = 5\%$$

$$1-\alpha = \text{level of confidence} \\ = 95\%$$

NOW,

95% confidence interval for π is given by,

$$P \pm Z_c \sqrt{\frac{pq}{n}} \\ = 0.7 \pm 1.96 \sqrt{\frac{0.7 \times 0.3}{500}} \\ = 0.7 \pm 0.0402$$

$$= 0.7402, 0.6598$$

$$= 74.02\%, 65.98\%$$

Conclusion:

There are between 65.98% to 74.02% of people who visit the mall make payment through credit card, and we are 95% sure in saying so.

Sample size determination ($n = ?$)

→ It is important to decide the no. of items to be selected from the population to constitute the sample. This is because the accuracy of estimation is directly related with sample size.

The size of the sample must be optimum, considering accuracy, representativeness and reliability.

Factors affecting the sample size:

1. Variability of characteristic in the population
2. Size of the population
3. Margin of error
4. Confidence probability
5. Sampling design
6. Nature of analysis
7. Number of categories or classes.

Sample size determination for studying mean (μ)

1. σ known case

The $(1-\alpha) \times 100\%$ confidence interval for popn mean (μ) is given by,

$$\bar{x} \pm Z_c \frac{\sigma}{\sqrt{n}}$$

$$\therefore E = Z_c \frac{\sigma}{\sqrt{n}}$$

$$\text{or } \sqrt{n} = Z_c \frac{\sigma}{E} \quad \therefore n = \frac{Z_c^2 \cdot \sigma^2}{E^2}$$

| E = Desired margin of error.

If N is known we used second approximation

$$n_0 = \frac{n}{1 + n/N}$$

2) σ unknown case

The $(1 - \alpha) \times 100\%$ confidence interval for μ is given by,

$$\bar{x} \pm t_{\alpha/2} (n-1) \frac{s}{\sqrt{n}}$$

E = Margin of error

$$= t_c \cdot \frac{s}{\sqrt{n}}$$

$$\text{or, } \sqrt{n} = \frac{t_c \cdot s}{E}$$

$$\therefore n = \frac{t_c^2 \cdot s^2}{E^2}$$

If N is known, we use

$$n_0 = \frac{n}{1 + n/N}$$

Sample size determination of studying mean (μ)

Numerical Problem.

- # A cable TV company would like to estimate the average no. of hours its customers spend watching TV per day. What size sample is needed if the company wants to have 95% confidence that its estimate is correct to within ± 15 minutes. The S.D. estimated from previous studies is 3 hrs.

→ Solution:

x = No. of hrs of customers spend watching TV per day.

n = Desired sample size for study = ?

6 = Prior estimate of population S.D.
= 3 hrs.

E = Desired margin of error
= ± 15 minutes.

α = Level of significance = 5%.

$1 - \alpha$ = Confidence probability = 95%.

Z_c = Critical value of z for $\alpha = 5\%$.
= 1.96

Thus the desired sample size is given by,

$$n = \frac{Z_c^2 \cdot 6^2}{E^2}$$

$$= (1.96)^2 \times (180)^2 \\ (15)^2$$

$$= 553.1904 \approx 554$$

∴ The study must include 554 people or viewers in order to estimate mean viewing time.

Note:

1. If N is known, find estimate of n , $n_0 = \frac{n}{1 + n/N}$
2. If σ is known use t instead of z .

Sample size determination for studying population proportion of success (π)

→ The $(1 - \alpha) \times 100\%$ confidence interval for population proportion of success (π) is given by,

$$p \pm z_c \sqrt{\frac{pq}{n}}$$

The margin of error is given by.

$$E = z_c \sqrt{\frac{pq}{n}}$$

$$\text{or, } \sqrt{n} = z_c \sqrt{\frac{pq}{E^2}}$$

$$\therefore n = \frac{z_c^2 \times pq}{E^2}$$

Note:

1. If N is known, find estimate of n , $n_0 = \frac{n}{1 + n/N}$
2. If there is no prior estimate of π , use $p=0.5$, $q=0.5$ in the formula.

Numerical

A survey is planned to study customer's satisfaction concerning after-sale support and service of product. A pilot study was conducted (with small sample size) and found out that an initial estimate of % of customer not satisfying with after sale support and service is 38%. What sample size is needed to conduct new study for proportion of customers satisfied and to obtain a margin of error of 4% and at 95% confidence level?

→ Solution:

x = customer satisfaction (satisfied, not-satisfied)
success - NOT

n = No. of customers required for sample survey = ?

$$p = \text{prior estimate of population of success} \\ = 38\% = 0.38$$

$$q = \text{prior estimate of population of failure} \\ = 1 - p \% = 1 - 0.38 = 0.62$$

$$\text{Desired margin of error for the study} = \pm 4\% \\ = 0.04$$

$$\text{confidence probability } (1-\alpha) = 95\%.$$

$$\text{significance probability } (\alpha) = 5\%.$$

$$Z_c = \text{critical value of } Z \text{ for } \alpha = 5\% \\ = 1.96.$$

∴ Required sample size is given by,

$$n = \frac{Z_c^2 \times p \times q}{E^2} \\ = \frac{(1.96)^2 \times 0.38 \times 0.62}{(0.04)^2} \\ = 565.6756 \approx 566.$$

conclusion:

The new study must include 566 people in order to estimate proportion of customer who are not satisfied with the service of the company.

UNIT-2

Hypothesis Testing

A hypothesis may be defined simply as a statement about one or more populations. It's the prediction of outcome of research.

The purpose of hypothesis testing is to aid the researcher, administrator etc in reaching a conclusion concerning a population by examining a sample from that population.

Hypothesis is a hunch, guess, imaginative idea that arise as a result of a prior thinking about the subject or topic, examination of the available data and material including related studies and discussion with counsel of experts and interested parties, which forms the basis for action or investigation and validity of which remains to be tested.

A researcher is concerned with two types of Hypotheses.

- Research hypothesis is the conjective or supposition that motivates the research.
- Statistical hypothesis is statement or claim or assertion about pop" parameter that may be evaluated by appropriate statistical techniques.

converting research hypothesis into a statistical hypothesis.

Research hypothesis

↓
what variable (response/outcome measurement) can best test the Hypothesis?

↓
which summary statistic (mean, median, proportion or correlation, etc) should be used? This statistic depends on the response / outcome variable.

↓
Re-phrase the hypothesis into statistical terms

Hypothesis is stated in terms:

1. Null Hypothesis (H_0, H_N):

It is the hypothesis to be tested. It is designated by the symbol H_0 . It is sometimes referred to as 'hypothesis of no difference', since it is a statement of agreement with (or no difference from) conditions presumed to be true in the population of interest. It nullifies the claim that experimental result is different from the one observed already.

2. Alternative Hypothesis (H_1, H_A):

The alternative hypothesis is set up as the opposite of the null hypothesis and represents the conclusion supported if the null hypothesis is rejected. It is designated by H_1 . If sample information fail to support the hypothesis, then we must conclude that something else is true. Whenever we reject the test hypothesis, the conclusion we do accept is called the alternative

hypothesis.

Example 1: Specific value testing hypothesis.

Research Hypothesis - After taking a particular medication the heart beat will increase.

Statistical Hypotheses -

$H_0 : \mu = 72$ (No change in heart beat)

$H_1 : \mu \neq 72$ (Two-tailed test / Two-sided test)

or, $H_1 : \mu < 72$ (Left-tailed test / Left-sided test)

or, $H_1 : \mu > 72$ (Right-tailed test / Right-sided test)

Example 2: Comparative Hypotheses.

Research Hypotheses - Vianet is better than worldlink.

variable

↓ →

categorical (Satisfied / Not satisfied)	Numerical (Band width at particular time)
--	--

$H_0 : \pi_1 = \pi_2$

π_1 = proportion of satisfied customer of vianet

π_2 = proportion of satisfied customer of worldlink

$H_1 : \pi_1 \neq \pi_2$

or, $H_1 : \pi_1 < \pi_2$

or, $H_1 : \pi_1 > \pi_2$

$H_0 : \mu_1 = \mu_2$

μ_1 = Average bandwidth of vianet by using specific package

μ_2 = Average bandwidth of worldlink

$H_1 : \mu_1 \neq \mu_2$

or, $H_1 : \mu_1 < \mu_2$

or, $H_1 : \mu_1 > \mu_2$

Example 3 : Relational Hypotheses

Research Hypothesis - Prolonged use of computer will produce eye problem

Statistical Hypothesis -

H_0 : Eye problem is independent of duration use of computer.

H_1 : Eye problem is not independent (dependent) of duration use of computer.

Methods of Testing Hypotheses.

Parametric

Critical value approach p-value approach
(Modern approach)

Non-parametric Methods

Critical value approach p-value approach

Steps in testing hypotheses (critical value approach)

Step 1 : State the null and alternative hypothesis

H_0 = Null Hypotheses

H_1 = Alternative Hypotheses.

Step 2 : choosing α level of the test

α = Prob { Type I error }

There is always risk associated with decision making using hypothesis testing method. We make two types of errors in hypothesis testing procedure.

Decision from sample	True state	
H_0 True	H_0 True $(1 - \alpha)$	H_0 False (H_1 True) (Type I error) (β)
H_0 False (H_1 True)	Correct Decision (Type II error) (α)	Wrong Decision $(1 - \beta)$

$$\alpha = \text{Prob} \{ \text{Type I error} \}$$

= Prob { Reject H_0 / H_0 True } $\quad \alpha \downarrow \beta \uparrow$
 decision from true state $\quad \alpha \uparrow \beta \downarrow$
 sample

$$\alpha = \text{Significant probability}$$

$$1 - \alpha = \text{Prob} \{ \text{Accept } H_0 / H_0 \text{ True} \}$$

$$\beta = \text{Prob} \{ \text{Type II error} \}$$

$$= \text{Prob} \{ \text{Accept } H_0 / H_0 \text{ False} \}$$

$$1 - \beta = \text{Power of the test}$$

$$= \text{Prob} \{ \text{Reject } H_0 / H_0 \text{ False} \}$$

Choosing α level of the test

→ Type I error is usually fixed in advance by choosing appropriate level of significance (α) employed for the test. There is no single standard level of significance for testing

hypotheses. It is possible to test hypothesis at any level of significance. Typically we choose $\alpha = 5\%$ or $\alpha = 1\%$. choosing β

we cannot directly control our risk of type II error. One way to reduce β 's is to increase sample size. Large sample generally reduces type II risk. The other way to control β risk is to increase α risk.

Step 3: choose appropriate test statistics

The test statistics is a certain formula for the particular test which serves as a decision maker. Since decision to reject or accept H_0 depends on the magnitude of it.

Step 4: finding the tabulated or critical value of test statistics.

The critical value or tabulated value divides the distribution in two types:

- I. Acceptance region
- II. Rejection region

Step 5: finding the actual value or computed value of test statistics

using the sample information, we compute a observed value of the test statistics.

Step 6: Statistical decision

We compare observed value of test statistic with critical value of test statistic and we decide whether to reject or not the hypothesis.

Step 7: Conclusion

In this step we make some inferences about the population of interest. (Knowledge discovery)

Types of parametric test

1. Z-test for mean of population, σ known
(single sample case)
2. t-test for population mean (μ), σ unknown
(single sample case)
3. Z-test for difference of two population means
(μ_1 and μ_2) σ_1 and σ_2 known
(Two independent samples)
4. t-test for difference of two population means
(μ_1 and μ_2) σ_1 and σ_2 unknown
(Two independent samples)
5. Paired t-test (t test for difference of two population means) (Two related / dependent samples)
6. Z-test for population (π), (single sample case)
7. Z-test for diff. of two population proportions
(π_1 and π_2) (two independent samples)

I. Z-test for mean of population, σ known (single sample case)



Function of the test

The function of the test is to determine whether a sample drawn from the population have a specific mean (μ_0), provided that pop's popn in SD (σ) is known.

Test assumptions

1. The variable of interest is normally distributed in the population.
2. Sample is drawn randomly from the population.
3. Measurement scale is atleast interval but ratio scale data is not.
4. Population SD is known

Hypotheses to test

Let μ be the mean of characteristic in the population from which a random sample is drawn.

- Null Hypothesis

$H_0: \mu = \mu_0$ (The popn has specific mean value μ_0)

- Alternative Hypothesis

$H_1: \mu \neq \mu_0$ (The popn doesn't have specific mean value μ_0)

OR, $H_1: \mu > \mu_0$ (Popn has mean value significantly higher than μ_0)

$H_1 \neq$ Two-sided test
 > Right ... "
 < Left ... "

Date _____
Page _____

OR, $H_1: \mu < \mu_0$ (Popn has mean value significantly lower than μ_0)

Test statistics

We know. $\bar{x} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$

Then, Z transformation \bar{x} is given by,

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \rightarrow \text{Test statistics}$$

\bar{x} = sample.

μ_0 = value of μ assumed under H_0 .

Z has standard normal deviation (SND) with mean 0 and SD of 1.

$$Z \sim \text{SND}(0, 1)$$

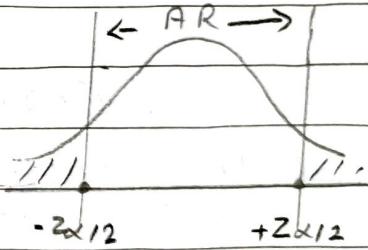
Decision Rule

1. $H_0: \mu = \mu_0$ (Two-sided test)

$$H_1: \mu \neq \mu_0$$

Reject H_0 if $| \text{cal } Z | \geq z_{\alpha/2}$

Reject H_0	Accept H_0	Reject H_0
$-z_{\alpha/2}$		$+z_{\alpha/2}$

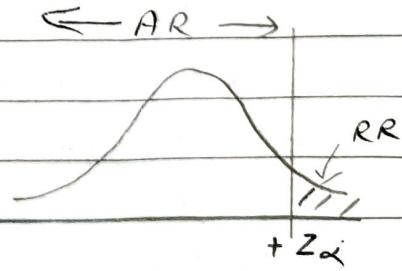


2. $H_0: \mu = \mu_0$

$$H_1: \mu > \mu_0$$

Reject H_0 if $\text{cal } Z > z_\alpha$

Accept	Reject
	$+z_\alpha$



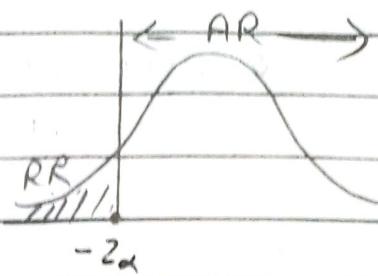
3. $H_0: \mu = \mu_0$

$H_1: \mu < \mu_0$

Reject H_0 if $\text{cal } Z \leq -Z_\alpha$

Reject H_0	Accept H_0
--------------	--------------

$-Z_\alpha$



Critical z

α	Two-sided test	Right-sided test	Left-sided test
3%			
5%	-1.96 +1.96	+1.65	-1.65
10%	-2.058 +2.058	+2.33	-2.33

Numerical

A manufacturer of a certain brand of 9 volt batteries claims that the average life of its battery is 40 hours when used in electronic devices with the SD of 5 hours. To test the manufacturer's claim, a random sample of 100 batteries was tested and it showed an average life of 38 hours. What can you conclude about manufacturer's claim at 5% level of significance?

→ Solution:

Step 1: Setting up Null and Alternative Hypotheses
Here,

$X = \text{life of 9 volt batteries}$

Let μ be the average life of 9-volt battery manufactured by the company.

So will it be different?
Not equal to or equal? } two tailed test

{ Will it be greater than?
Will it be higher than? } Right tailed test

{ Will it be lesser than?
Less than? } Left tailed test

Date _____
Page _____

Null Hypothesis

$H_0: \mu = 40$ (Batteries manufactured by company has specific mean life of 40 hrs)

Alternative Hypothesis

$H_1: \mu \neq 40$ (The mean life of batteries manufactured by company is not equal to 40)

Step 2: choice of α for the test

α = level of significance employed for the test
= Prob { Type I error }
= 5% (Given)

Test statistics (Step 3)

The appropriate test statistics for the test is

$$Z = \frac{\bar{X} - 40}{\sigma/\sqrt{n}}$$

Z has SND with mean 0 and SD of 1.

Step 4: Critical Z or Tabulated Z .

Here, $\alpha = 5\%$ (and test is two sided)

From the table of Z distribution, the critical Z is,

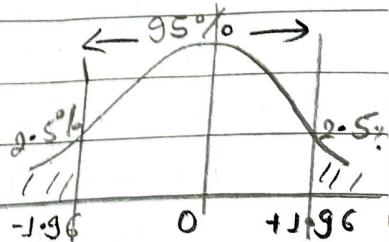
$$1.96$$

Reject H_0	Accept H_0	Reject H_0
-1.96		+1.96

$$AR = -1.96 < Z < +1.96$$

$$RR = Z \geq +1.96 \text{ OR } Z \leq -1.96$$

(CR)



Step 5: H₀ calculated Z

Now,

$$Z = \bar{x} - \mu_0$$

$$6/\sqrt{n}$$

$$= 38 - 40$$

$$6/\sqrt{100}$$

$$= -4$$

Step 6: Statistical Decision

Since cal Z = -4 falls in the lower rejection region ($Z \leq -1.96$) we reject H₀ at 5% level of significance in favour of H₁.

Step 7:

The manufacturer's claim is not valid. Then, mean life of 9 volt batteries is actually different from 40 hours.

2. t-test for population mean (μ), σ unknown
(single sample case)



Function of test

The function of the test is to determine whether a sample drawn from the population have a specific mean, provided that its population standard deviation σ is unknown.

Test assumptions

1. Variable of interest is normally distributed in the

population.

2. Sample is drawn randomly from the population.
3. Measurement scale is at least interval, but ratio scale data is preferred.
4. Population standard deviation σ is unknown.

Hypothesis to test

Let μ be the mean of the characteristics in the population from which a sample of n observations is drawn.

$H_0: \mu = \mu_0$ (population has specific mean value) against

$H_1: \mu \neq \mu_0$ (population do not have specific mean value μ_0)

$H_1: \mu < \mu_0$ (population have mean value significantly lower than μ_0)

$H_1: \mu > \mu_0$ (population have mean value significantly higher than μ_0)

Test statistics

We know when population standard deviation (σ) is not known, the sampling distribution of the sample mean \bar{x} follows students t-distribution with $n-1$ degrees of freedom. Hence, the appropriate test statistics for the test is given by,

$$t = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

where \bar{x} = sample mean

μ_0 = hypothesized value of mean

s = known population standard deviation

n = size of the sample drawn from the population

Distribution of test statistic

If null hypothesis is true the quantity $\frac{\bar{x} - \mu_0}{s/\sqrt{n}}$ follows

Student's t-distribution with $n-1$ degrees of freedom.

Decision Rule

we adopt following decision rule:

case I (Two sided test) : Reject H_0 if $|t| > t_{\alpha/2}(n-1)$

Case II (left sided test) : Reject H_0 if $t \leq -t_{\alpha}(n-1)$

Case III (Right sided test) : Reject H_0 if $t \geq +t_{\alpha}(n-1)$

Two tailed test

Reject H_0	Do not reject H_0	Reject H_0
$-t_{\alpha/2}(n-1)$		$+t_{\alpha/2}(n-1)$

left tailed test

Reject H_0	Accept H_0
$-t_{\alpha}(n-1)$	

Right tailed test

Accept H_0	Reject H_0
	$+t_{\alpha}(n-1)$

Numerical:

- # A random sample of soil specimens was obtained and the amount of organic matter (%) in the soil was determined for each specimen, resulting in the accompanying data.

1.10	5.09	0.97	1.59	4.60	0.32	0.55	1.45
0.14	4.47	1.20	3.50	5.02	4.67	5.22	2.69
3.98	3.17	3.03	2.21	0.69	4.07	3.81	1.17
0.76	1.17	1.57	2.62	1.66	2.05		

Does this data suggest that the true average percentage of organic matter in such soils is something other than 3%? Carry out a test of the appropriate hypothesis at significance level 0.01. Would your conclusion be different if 0.05 had been used? calculate p-value.

- # What is d.f.?

→ No. of observations we can freely specify to calculate a particular statistic is called its degree of freedom. Let statistics be \bar{x} and $n=3$.
 $x_1, x_2, x_3 \quad \bar{x} = 80$

$$\text{Let } x_1 = 15, x_2 = 25, x_3 = ?$$

$$\text{Restriction} \rightarrow \sum x = n \cdot \bar{x}$$

$$x_1 + x_2 + x_3 = 3 \times 80$$

$$15 + 25 + x_3 = 240$$

$$x_3 = 200$$

$$\bar{x} = \sum x / n$$

$$\sum x = n \cdot \bar{x}$$

General \bar{x} has

$n-1$ degrees of freedom.

General formula of degree of freedom.

(P) = sample size - no. of restrictions

→ SOLUTION:

Step 1: After setting up null and alternative hypothesis.

X = Amount of organic matter in the soil (%)

Let μ be the mean amount of organic matter in the soil (%)

Null Hypothesis

$H_0: \mu = 3\%$ (Mean organic matter is 3%)

Alternative Hypothesis

$H_1: \mu \neq 3\%$ (Mean organic matter is other than 3%)

Step 2: choice of α for the test. If α is not given

α = level of significance

then test H_0

= prob{Type II error} at $\alpha = 5\%$

= 0.01

Step 3: Test statistic

The appropriate test statistics is given by,

$$t = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

$$S/\sqrt{n}$$

where \bar{X} =

$$\mu_0 =$$

$$S =$$

$$n =$$

The test statistic follows student's t-distribution with $n-1$ df.

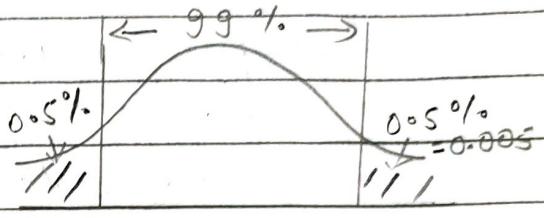
Step 4: Critical value of t
(tabulated value of t)

Hence, $\alpha = 1\% = 0.01$ and the test is two-sided. Hence, there are two rejection regions.

$$df = n - 1$$

$$= 30 - 1 = 29.$$

Reject H_0	Accept	Reject H_a
-2.756	+2.756	



$$AR: -2.756 < t < +2.756$$

$$RR: t \geq +2.756 \text{ OR } t \leq -2.756$$

Step 5: Calculation of t
(observed value of t)

Hence

$$\sum X = 74.44 \quad n = 30$$

$$\sum X^2 = 260.409$$

$$\text{Now, } \bar{X} = \sum X / n = 74.44 / 30 \\ = 2.4813$$

$$S = \sqrt{\frac{1}{n-1} \left\{ \sum X^2 - n \cdot \bar{X}^2 \right\}}$$

$$= \sqrt{\frac{1}{30-1} \left\{ 260.409 - 30 \times 2.4813^2 \right\}}$$

$$= 1.6156$$

$$\therefore t = \frac{\bar{X} - \mu_0}{S / \sqrt{n}} = \frac{2.4813 - 3}{1.6156 / \sqrt{30}} \\ = -1.76$$

Step 6: Statistical Decision.

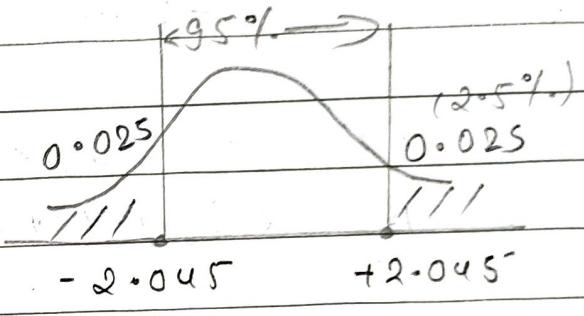
Since calc $t = -1.76$ falls in AR $(-2.756 < t < +2.756)$, we do not reject H_0 at 1% level of significance.

Step 7: Conclusion

The mean organic matter is not different from 3%.

$$(b) \alpha = 5\%$$

Reject	Accept	Reject
-2.045	+2.045	
-1.76		



$$\alpha(1) \rightarrow$$

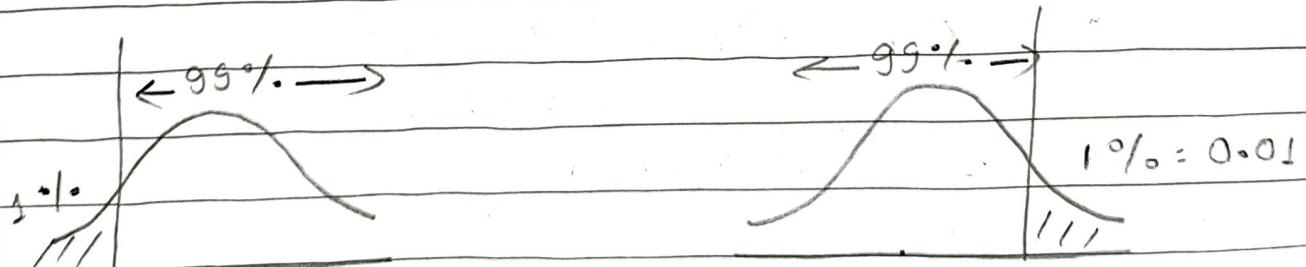
$$t = -1.76$$

$$H_0: \mu = 3\%$$

$$H_1: \mu < 3\%.$$

$$H_0: \mu = 8\%$$

$$H_1: \mu > 8\%.$$



Reject	Accept
-2.462	1.76

$$\alpha(1) \rightarrow 0.05$$

Accept	Reject
-1.76	+2.462

$$\alpha(1) \rightarrow 0.01$$

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\text{SP} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\text{SP} = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}$$



3. Two sample z test for different mean.

Function of the test

The function of the test is to determine whether two samples drawn from the two independent populations have same mean or not provided that population SDs are known.

Test assumptions.

1. Sampling method for each sample is simple random sampling.
2. Samples are independent.
3. Two populations are normally distributed if samples are of sizes less than 30 but do not require this condition if sample are larger than 30 (The sampling distribution of diff of two means is normally distributed, generally it is true by virtue of central).
4. Population SDs are known.

Hypothesis to test.

Let μ_1 and μ_2 be the population mean of the population I and population II respectively from which samples are drawn independently.

Two tailed test	Left tailed test	Right tailed test
$H_0: \mu_1 = \mu_2$	$H_0: \mu_1 \geq \mu_2$	$H_0: \mu_1 \leq \mu_2$
$H_1: \mu_1 \neq \mu_2$	$H_1: \mu_1 < \mu_2$	$H_1: \mu_1 > \mu_2$

Test statistic

The appropriate test statistic is given by,
 $Z = \frac{\text{Sample diff. of means} - \text{Hypothesized diff. of mean}}{\text{Standard error of diff. of means}}$

$$= \frac{\bar{x}_1 - \bar{x}_2 - \mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$= \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

In practice, the two sample z-test is not often used because the two popn standard deviations σ_1 and σ_2 are usually unknown. Instead, sample SDs are and the t-distribution are used. However, if samples are of large size, even though popn SDs are not known, we can still use z-distribution replacing popn SDs by respective sample SDs.

The test statistic is given by:

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Numerical :

To compare the starting salaries of college graduates majoring in engineering and computer science a random sample of recent college graduates in each major were selected and the following info obtained.

Major	Mean(\$)	SD(\$)
Engineering	56,202	2225
compt'r science.	50,657	2375

Do the data provide sufficient evidence to indicate a diff' in average starting annual salaries for college graduates who majored in engineering and computer science? Test using $\alpha = 0.05$.

→ Solution :

Step 1 : setting up Null and Alternative Hypothesis.
 Let X_1 = starting annual salary of engineering major.

X_2 = starting annual salary of computer science major.

Let μ_1 and μ_2 be the mean starting annual salaries of engineering major and computer science major respectively.

Null Hypothesis

$H_0: \mu_1 = \mu_2$ (There is no significant difference between starting annual salaries of engineering major and computer science major)

Alternative Hypothesis

$H_1: \mu_1 \neq \mu_2$ (There is significant difference)

Step 2: choice of α for the test

$$\begin{aligned}\alpha &= \text{level of significance of the test} \\ &= \text{Prob}\{\text{Type I error}\} \\ &= 0.05\end{aligned}$$

Step 3: Test statistics

The appropriate test statistic for the test is given by,

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

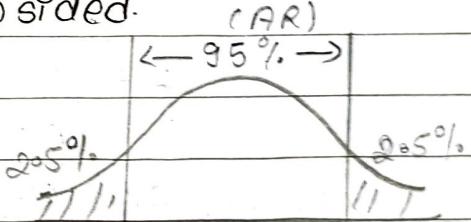
The test statistic has SND with mean = 0 and SD = 1.

Step 4: Critical z (Tabulated z)

Here, $\alpha = 5\%$ and test is two-sided.

$$AR: -1.96 < Z < +1.96$$

$$\begin{aligned}RR: Z &\geq +1.96 \text{ OR} \\ &Z \leq -1.96\end{aligned}$$



Reject H ₀	Accept H ₀	Reject H ₀
-1.96		+1.96

Step 5: Calculated z (observed z)

$$\begin{aligned}Z &= \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \\ &= \frac{56202 - 50657}{\sqrt{\frac{(2255)^2}{50} + \frac{(2375)^2}{50}}}\end{aligned}$$

$$= +11.97$$

Step 6 : Statistical Decision

Since cal Z (11.97) is very higher and falls in the upper critical region ($2.3 + 1.96$), we strongly reject H_0 in favour of H_1 .

Step 7 : Conclusion

There is significant diff. in the mean starting annual salaries of engineering major and computer science major. Hence engineering majors are paid more than computer science major.

4. Independent sample t-test



Function of the test

The independent t-test compares the means of two or more independent groups in order to determine whether there is statistical evidence that the associated population means are significantly different.

The independent samples t-test is a parametric test.

Common use:

The independent samples t-test is commonly used to test the following:

- statistical differences between the means of two groups.
- statistical differences between the means of two interventions.

-statistical differences between the means of two change scores.

Data requirement / Test assumptions

1. Dependent variable that is continuous (e.g. interval or ratio level)
2. Independent variable that is categorical (two groups)
3. Independent samples / groups (e.g. independence of observations).
4. Random samples of data from the population
5. Normal distribution (approximately) of the dependent variable for each group
6. Homogeneity of variances (e.g. variances approximately equal across groups)
7. NO outliers.

Hypothesis to test

The null and alternative hypothesis are given below:

<u>Null hypothesis</u>	<u>Alternative hyp.</u>	<u>No. of tails</u>
$H_0: \mu_1 = \mu_2$ or, $H_0: \mu_1 - \mu_2 = 0$	$H_1: \mu_1 \neq \mu_2$	Two
$H_0: \mu_1 > \mu_2$ or, $H_0: \mu_1 - \mu_2 \geq 0$	$H_1: \mu_1 < \mu_2$	One (left tailed test)
$H_0: \mu_1 \leq \mu_2$ or, $H_0: \mu_1 - \mu_2 \leq 0$	$H_1: \mu_1 > \mu_2$	One (right tailed)

Test statistics

The test statistic for this test is given by,

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p} = \frac{\bar{x}_1 - \bar{x}_2}{s_p}$$

$$s_p = \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \sqrt{\frac{s_1^2 + s_2^2}{n_1 + n_2}}$$

where,

$$s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$$

\bar{x}_1 = Mean of first sample

\bar{x}_2 = Mean of second sample

n_1 = Sample size (i.e. numbers of observations) of first sample

n_2 = sample size (i.e. no. of observations) of second sample.

s_1 = S.D. of first sample.

s_2 = S.D. of second sample.

s_p = Pooled SD. i.e. pooled estimate of common SD

6.

Decision Rule

Hypothesis

Case I (Two-sided)

Reject H_0 if $|cal{t}| \geq t_{\alpha/2}(n-1)$

case II (left ")

Reject H_0 if $cal{t} \leq -t_{\alpha}(n-1)$

case III (Right ")

Reject H_0 if $cal{t} \geq t_{\alpha}(n-1)$

Large sample case: If sample sizes are large, even though population S.D.s s_1^2 and s_2^2 are unknown, t-test can be replaced by z-test using following formula.

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Example

- # In a packaging plant, a machine packs carton with jars. It is supposed that a new machine will pack faster on the average than the machine currently used. To test that hypothesis the time it takes each machine to pack ten cartons are recorded. The results seconds, are shown in the following table:

New machine	Old machine
42.1, 41.3, 42.4, 43.2,	42.7, 43.8, 42.5, 43.1,
41.8, 41.0, 41.8, 42.8,	44.0, 43.6, 43.3, 43.5,
42.8, 42.7	41.7, 44.1

Do the data provide sufficient evidence to conclude that, on the average, the new machine packs faster? Perform the required hypothesis test at the 5% level of significance.