

Unit - 3

Sampling :-

- Fundamental concept in statistics that involves selecting subset of individuals / items from a large population.
- It is used to draw conclusion about entire population.

Population - The entire group of individuals, items or data points under consideration.

Sample - A subset of population for analysis.

Types of Sampling

Probability Sampling

A probability sample means that every member of the population has a chance of being selected.

- used in quantitative research

Four types of probability sampling -

1 Simple random sampling

2 Systematic sampling

3 Stratified sampling

4 Cluster sampling

Non Probability Sampling

Individuals are selected on non-random criteria & not every individual has a chance of being selected.

- easier & cheaper to access

- high risk of sample biasness

- used in exploratory &

qualitative research

Four types of Non-probability Sampling -

1) Convenience Sampling

2) Purposive Sampling

3) Snowball Sampling

4) Quota Sampling

Probability Sampling techniques -

1. Simple random sampling :- Every member of the population has equal chance

- sample frame includes whole population.

- Use tools like random number generator

eg There are 1000 employees & selecting 100 out of these 1000.

2. Systematic Sampling :- Similar to simple random sampling

- selection is done by choosing individual at regular interval.

eg Out of 100, 10 students are chosen such as (6th, 16th, 26th, ..., 86th, 96th) student is selected in sample.

3. Stratified Sampling :- Dividing the population into subpopulation based on relevant criteria.

- Then sample is created from subgroup (subpopulation) using systematic sampling.

eg Out of 1000, 800-female & 200-male then sample of 100 is generated 80-female & 20-male. (Strata based on gender).

4 cluster sampling :- Dividing the population into subgroups, but each subgroup must have similar characteristics

- Instead of randomly selecting individuals, we select entire subgroups as sample.

eg company has offices in 10 cities, you can make sample by selecting 3 offices (3 cities). (Here every office is a cluster).

* Multistage sampling - If the cluster is too much larger, then sample individuals from within each clusters.

this is known as Multistage sampling.

Non-Probability Sampling -

2 Convenience Sampling :- Sample simply includes the individuals who happen to be most accessible to the researcher.

- Easy & inexpensive way to gather initial data. (There is no way

to tell if sample represent population)

eg College surveys like asking nearby students about certain things regarding college.

2. Voluntary

2. Purposive Sampling - Also known as judgemental sampling, involves the researcher using their expertise to select a sample that is most useful to purpose of the research.

- Also used in qualitative research, where the researcher wants to gain detailed knowledge about a specific phenomenon. (population is small & specific).

eg Taking opinion of 4th year students regarding placement

3. Snowball sampling - Used when population is hard to access. Recruiting participants via other participants. Since people increase one after another like "snowball", your access increases.

- Sampling biasness happens in this.

eg Gathering data of homeless people from one homeless person, thus the sample increases.

4. Quota Sampling - Relies on the non-random selection of a predetermined number or proportion of units. This is called quota.

- First divide population into mutually exclusive subgroups (strata) & then sample units until you reach your quota.

eg population - 1000 divide stratas - vegetarian, non-veg. & vegans recruit 200 for each then compare.

Difference b/w Parameters & Statistic

Parameter

Statistic

- Numerical value that describe characteristics of population
- Numerical value that describe characteristics of sample.
- Represented by Greek letters
- Represented by Roman letters
- mean = μ s.d. = σ size = N Mean = \bar{x} s.d. = s size = n
- Parameters are fixed & unchangeable, usually known
- Statistics can be variable depending on different samples
- (calculated using data from entire population)
- Coupled for different samples acc. to data.
- Used to make inference about a population
- Used to make inference about population based on a sample.

Test of Significance

- Statistical method used to determine whether an observed effect (result) is likely to be genuine or occurred by chance.
- These test help to draw meaningful conclusion from data.
- Null Hypothesis (H_0): Hypothesis is correct but rejected
- Alternate Hypothesis (H_1): Hypothesis is wrong but accepted.

- Type I error: (α) Rejecting H_0 while H_0 is true.

Null Hypothesis - Hypothesis which is tested for possible rejection under the assumption that it is true, is called type I error.
Null hypothesis.

- Type 2 error (β): Accepting H_0 while H_0 is false.

Alternate Hypothesis - Hypothesis which is tested for possible acceptance under the assumption that it is false, is called type 2 error.
Alternate hypothesis.

Type 1 error (α): The error of rejecting H_0 when H_0 is true is called type 1 error.

- It gives false positive conclusion.
- It is denoted by alpha (α).

$$\alpha = P(\text{Rejecting } H_0 \mid H_0 \text{ is true})$$

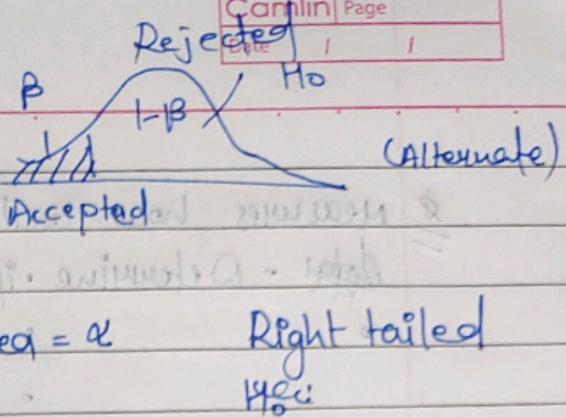
Type 2 error (β): The error of accepting H_0 when H_0 is false is called type 2 error.

- It gives false negative conclusion.
- It is denoted by beta (β).

$$\beta = P(\text{Accepting } H_0 \mid H_0 \text{ is false})$$

Critical region: If the statistic value falls in critical region, the null hypothesis is rejected.

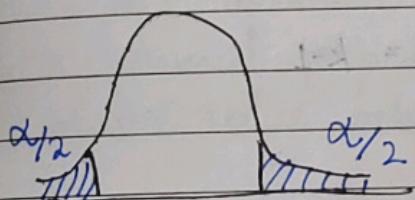
Critical value: Boundary value that separates critical region from acceptance region.



$$\text{Area} = \alpha$$

Acceptance = $1-\alpha$ \rightarrow left tailed

Two tailed



$$\text{Area} = \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha$$

$$H_0: \mu = \bar{x}$$

$$H_1: \mu \neq \bar{x}$$

- Significance level: The probability of rejecting the null hypothesis when it is true. (Standard val)
 - standard value is 0.05 or 5%.
- P-value: The observed probability of rejecting the null hypothesis when it is true. (Experimental val)
 - If P-value < Significant value \rightarrow Null hypothesis rejected
 - If P-value \geq Significant value \rightarrow Null hypothesis accepted

Goodness of fit test

A goodness of fit test is a statistical method used to determine if an observed frequency distribution fits a theoretical model or expected distribution.

- Sample data follows a specific probability distribution.

- ❖ Measures how well observed data matches a fitted model.
- Determine if a sample follows a normal distribution.
 - If categorical variables are related
 - Random samples are from same distribution
- Degrees of Freedom: The goodness of fit test depends on the number of categories & constraints in the problem.

$$\text{degrees of freedom (df)} = k-1$$

❖ Chi Square test (Pearson's chi sq. test) χ^2

- Non parametric test
- Pearson's chi-sq. test is a statistical test for categorical data. It is used to determine whether your data are significantly different from what you expected.

There are two types of Pearson's chi-sq. test

- 1) Chi-sq. goodness of fit test
- 2) Chi-sq. Test of Independence

25. i) Chi-sq. Goodness of fit test

- It is used to test whether the frequency distribution of a categorical variable is different from your expectation.

ii) Chi-sq. Test of Independence

- It is used to test whether two categorical variables are related to each other.

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} \quad df = n-1$$

where
 O_i = Observed frequency
 E_i = Expected frequency
 n = number of objects

*# T-test ($n < 30$)

To test whether the mean of a sample drawn from a normal population deviates significantly from a stated value when variance of population is unknown

One Sample

$$t = \frac{|\bar{x} - \mu|}{S/\sqrt{n}}$$

$$df = n-1$$

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

$$df = n_1 + n_2 - 2$$

\bar{x} = mean of sample \bar{x}_1, \bar{x}_2 = mean of sam₁, sam₂

μ = mean of Population

s_1, s_2 = sd of sam₁, sam₂

s = sd. of sample

n_1, n_2 = sam₁, sam₂ size

n = sample size

*# Z-test ($n > 30$)

- To test whether there is significance difference b/w sample mean & population mean, given that variance of population is known.

$$z = \frac{|\bar{x} - \mu|}{SEM}$$

\bar{x} = sample mean

μ = population mean

$$SEM = \frac{\sigma}{\sqrt{n}}$$

n = sample size. SEM = st. error of mean σ = sd. population

⇒ Two Sample Z-test

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

and because \bar{x}_1, \bar{x}_2 = sample means
 n_1, n_2 = sample size

Significance level with standard value for Z-test

Signi level	confi level	2-tailed	1-tailed
0.1	10%	90%	1.65
0.05	5%	95%	1.96
0.01	1%	99%	2.58
0.001	0.1%	99.9%	3.29

★ (5marks)

Ques. Write down the procedure of Hypothesis testing?

Ans.

Steps :-

- 1) Setting up Hypothesis —
 - Null Hypothesis
 - Alternate Hypothesis
- 2) Set up Significance level — (5%, 10%, 15%, ...)
- 3) Testing the hypothesis —
 - Z-test ($n > 30$)
 - t-test ($n < 30$) (or we have two means μ_1, μ_2)
 - χ^2 -test
 - F-test
 - ANOVA
- 4) Doing Computation (Comparing with Standard value)
- 5) Making Decision (Acceptance / Rejecting Null)

Ques. The mean life time of 81 glucometers produced by a company is found to be 58 months with a s.d. of 10 months. Company claims that mean life time of glucometers is 60 months. Test the hypothesis at 5% level of significance whether the 58 months life time is accepted or not.

Soln given $n=81$ $\mu=60$ $\bar{x}=58$
 $\sigma=10$ $\alpha=5\%$

H_0 : There is no significant difference b/w means $\bar{x}=\mu$
 H_1 : There is significant difference b/w means $\bar{x} \neq \mu$.

Test statistic $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$
 $= \frac{58 - 60}{10/\sqrt{81}} = \frac{-2}{10/9} = -1.8$

The standard value for two tailed test at 5% level of significance is 1.96
 $\text{since, } Z = -1.8 < 1.96$
 H_0 is accepted, means there is no significant difference b/w sample & population mean.

Ques A teacher claims that the mean score of student in his class is greater than 82 with a sd. of 20. If a sample of 81 students was selected at random with a mean score of 90 then check if there is enough evidence to support the claim at a significant level of 0.05.

Given - $n=81$ $H>82$ $\alpha=0.05$
 $\sigma=20$ $\bar{x}=90$

$\bar{x}_1 = 82$

H_0 : Mean score is greater than 82 $\mu > 82$

H_1 : Mean score is less than 82 $\mu < 82$

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

$$Z = \frac{|90 - 82|}{\frac{20}{\sqrt{18}}} = \frac{8}{\frac{20}{\sqrt{18}}} = 3.60$$

for one-tailed at 0.05 significant level
the standard value is 1.64

since, $Z = 3.6 > 1.64$

thus H_0 is rejected

Ques. Is there a significant difference in test score b/w 25 students who received inperson instruction & 25 persons who received online instruction. The mean test score for inperson group is 80 & for online group is 75 & the s.d. for inperson group is 5 & for online group is 7.

Given given, $\bar{x}_1 = 80$, $s_1 = 5$ and $n_1 = 25$

also given, $\bar{x}_2 = 75$, $s_2 = 7$, $n_2 = 25$

H_0 : There is no significant diff. b/w sample means $\bar{x}_1 = \bar{x}_2$

H_1 : There is significant diff. b/w sample means.

distribution for $t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$ where standard deviation of difference of sample means is $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

$$t = \frac{80-75}{\sqrt{\frac{25}{25} + \frac{49}{25}}} = \frac{5}{\sqrt{74}}$$

$$P(t > 2.91) = 0.01$$

$$P(t > 2.91) = 0.01$$

$$\text{since, } df = n_1 + n_2 - 2 = 25 + 25 - 2 = 48$$

$$sf = 0.05 \text{ standard value} = 2.01$$

$$\therefore t = 2.91 > 2.01$$

$$P(2.91 > 2.01) = 0.01$$

H_0 is rejected, that means there is significant diff. b/w means of two given samples.

A normal population has a mean of 6.8 & sd 1.5
a sample of 400 members gave a mean of 6.75, is there any significant difference.

$$\text{given } \mu = 6.8 \quad \bar{x} = 6.75 \quad sf = 0.05 \\ \sigma = 1.5 \quad n = 400 \quad H_0: \mu = \bar{x} \quad H_1: \mu \neq \bar{x}$$

$$z = \frac{|\bar{x} - \mu|}{\sigma/\sqrt{n}} = \frac{|6.75 - 6.8|}{1.5/\sqrt{400}} = \frac{0.05}{1.5/20} = 0.67$$

for two tailed at $sf = 0.05$ standard value = 1.96

$|z| = 0.67 < 1.96$ H_0 is accepted. There is no significant diff.

Que. The following table gives the number of accidents that took place in an industry during various days of week. Test if accidents are uniformly distributed over the week.

Days	(observed) No. of Accidents	Expected frequency	$(O-E)^2$	$\frac{(O-E)^2}{E}$
Monday	15	12	9	0.75
Tuesday	11	12	1	0.083
Wednesday	10	12	4	0.33
Thursday	12	12	0	0
Friday	10	12	4	0.33
Saturday	14	12	4	0.33
Sunday	13	12	1	0.083
	<u>84</u>			<u>2.326</u>

Soln

$$E = \frac{84}{7} = 12$$

H₀: Uniformly distributed
H₁: Non uniformly distributed.

$$\chi^2 = \sum \left\{ \frac{(O_i - E_i)^2}{E_i} \right\}$$

$$\chi^2 = 2.326$$

$$df = N-1 = 7-1 = 6$$

$$\text{standard value} = 12.59$$

since $\chi^2 = 2.326 < 12.59$ \therefore Null hypothesis accepted

Ques.

Day	Mon	Tue	Wed	Thus	Fri	Sat
Accidents	14	18	12	11	15	14
Expected	14	14	14	14	14	14
$(O_i - E_i)^2$	0	16	4	9	1	0
$\frac{(O_i - E_i)^2}{E}$	0	1.14	0.28	0.64	0.07	0

$$\chi^2 = \sum \left(\frac{(o_i - E_i)^2}{E} \right) = 2.13 < 11.07$$

$df = 6-1=5$ standard value = 11.07
 $\therefore H_0$ is accepted.

The life time of Electric bulbs for a random sample of 10 from a large consignment gave the following data can we accept the hypothesis that the avg life of bulb is 4000 hrs.

(in 1000)

Items	Life time	$n - \bar{x}$	Mean = $\frac{44}{10}$
i ₁	4.2	-0.2	0.04
i ₂	4.6	0.2	0.04
i ₃	3.9	0.5	0.25
i ₄	4.1	-0.3	0.09
i ₅	5.2	0.8	0.64
i ₆	3.8	-0.6	0.36
i ₇	3.9	-0.5	0.25
i ₈	4.3	-0.1	0.01
i ₉	4.4	0	0
i ₁₀	5.6	1.2	1.44
<u>44</u>			$sd = \sqrt{\frac{3.13}{9}}$

H_0 : There is no diff. b/w given & observed value $H_0 = 4000$

H_1 : There is diff.

$$\bar{x} = 4.4 \quad s = 0.58 \quad n = 10 \quad df = n-1 = 10-1 = 9$$

$$t = \frac{|4.4 - 4|}{\frac{0.58}{\sqrt{10}}} = \frac{0.4 \times \sqrt{10}}{0.58} = 2.18 < 2.262$$

thus null hypothesis is accepted.

Anova - Analysis of variance.

- When we have more than 2 samples

Ques. To assess the significance of possible variation in a certain test b/w the schools of a city a common test was given to a number of student taken at random from the 12th class of the three schools concerned. Make the analysis of variance from the below given results.

	A	B	C	M.G.0	S.S.	S.P.
1	2	3	4	20.0	2.0	0.8
2	4	5	6	18.0	2.0	1.0
3	6	7	8	16.0	2.0	0.8
				54.0	6.0	2.4

Null H₀: There is no significance difference $H_1 = H_2 = H_3$
Alternative H₁: There is significance difference

② calculate the mean of each sample

	A	B	C	$\bar{x}_A = \frac{12}{3} = 4$	$\bar{x}_B = \frac{15}{3} = 5$	$\bar{x}_C = \frac{18}{3} = 6$
1	2	3	4			
2	4	5	6	20.0	2.0	0.8
3	6	7	8	18.0	2.0	1.0
Total	12	15	18	54.0	6.0	2.4
Mean	4	5	6			

Calculate the ground average of means

$$\bar{X} = \bar{x}_A + \bar{x}_B + \bar{x}_C$$

$$\bar{X} = \frac{4+5+6}{3}$$

$$\bar{X} = 5$$

Calculation of SSC (sum of squares b/w samples)

$\bar{x}_A - \bar{X}$	$(\bar{x}_A - \bar{X})^2$	$(\bar{x}_B - \bar{X})$	$(\bar{x}_B - \bar{X})^2$	$(\bar{x}_C - \bar{X})$	$(\bar{x}_C - \bar{X})^2$
$4-5=-1$	1	$5-5=0$	0	$6-5=1$	1
$4-5=1$	1	$5-5=0$	0	$6-5=1$	1
$4-5=-1$	1	$5-5=0$	0	$6-5=1$	1

$$\sum (\bar{x} - \bar{\bar{X}}) = 3+0+3=6 \quad (\text{SSC})$$

↓ ground avg

Calculation of SSE (sum of squares within samples)

$A - \bar{x}_A$	$(A - \bar{x}_A)^2$	$(B - \bar{x}_B)$	$(B - \bar{x}_B)^2$	$(C - \bar{x}_C)$	$(C - \bar{x}_C)^2$
$2-4=-2$	4	$3-5=-2$	4	$4-6=-2$	4
$4-4=0$	0	$5-5=0$	0	$6-6=0$	0
$6-4=2$	4	$7-5=2$	4	$8-6=2$	4

$$\sum (x - \bar{x})^2$$

$$\sum (x - \bar{x})^2 = 8+8+8=24$$

Sq. of Variation	Sum of Squares	Degrees of freedom	Mean sum of squares	$F = \frac{MSC}{MSE}$
1. B/w Samples	SSC = 6	$V_1 = c-1$ $= 3-1 = 2$	$MSC = \frac{SSC}{c-1} = 3$	$F = 3/4$
2. Within Samples	SSR = 24	$V_2 = n-c$ $= 9-3 = 6$	$MSE = \frac{SSR}{n-c} = 4$	$= 0.75$

10. $n = \text{no. of samples} = 9$ $c = \text{no. of columns} = 3 (A, B, C)$
 since $0.75 < 5.14$, thus H_0 accepted

Anova table (c=1)

look	---	--- (E)	look column no. 2-p
sample no.	$n-c$	-----	-----
no.	-	(n-c) -	standard value.
-	-	(n-c)-1	

Ques. Calculate the Anova coefficient for following data

Plant	Number	Avg Span	SD
Hab	5	12	2
Masi	5	16	1
Rose	5	20	4

Given:

Plant	No.	Avg	SD
Hab	5	12 (\bar{x}_A)	2
Masi	5	16 (\bar{x}_B)	1
Rose	5	20 (\bar{x}_C)	4

$$\bar{\bar{x}} = \frac{\bar{x}_A + \bar{x}_B + \bar{x}_C}{3} = \frac{12 + 16 + 20}{3} = \frac{48}{3} = 16$$

$$\begin{array}{lll}
 (\bar{x}_A - \bar{\bar{x}}) = (12 - 16) = -4 & (\bar{x}_A - \bar{\bar{x}})^2 = 16 & \sum(\bar{x}_A - \bar{\bar{x}}) = 80 \\
 (\bar{x}_B - \bar{\bar{x}}) = (16 - 16) = 0 & (\bar{x}_B - \bar{\bar{x}})^2 = 0 & \sum(\bar{x}_B - \bar{\bar{x}}) = 0 \\
 (\bar{x}_C - \bar{\bar{x}}) = (20 - 16) = 4 & (\bar{x}_C - \bar{\bar{x}})^2 = 16 & \sum(\bar{x}_C - \bar{\bar{x}}) = 80
 \end{array}$$

since $n_1 = 5 = n_2 = n_3 \therefore \sum(\bar{x}_A - \bar{\bar{x}}) = (\bar{x}_A - \bar{\bar{x}}) \times n_1$

~~YTDN X+T+N : X+SSC + \sum(\bar{x} - \bar{\bar{x}}) = 80 + 0 + 80 = 160~~

Since $s^2 = \frac{\sum(x - \bar{x})^2}{n-1}$

~~for SSE we need $\sum(x - \bar{x})^2$~~

$$\therefore s^2(n-1) = \sum(x - \bar{x})^2$$

for $\sum(x_A - \bar{x}_A)^2 = s_1^2(n_1-1) = 4 \times (5-1) = 16$

$\sum(x_B - \bar{x}_B)^2 = s_2^2(n_2-1) = 1 \times (5-1) = 4$

$\sum(x_C - \bar{x}_C)^2 = s_3^2(n_3-1) = 16 \times (5-1) = 64$

$\therefore SSE = \sum(x - \bar{x})^2 = 16 + 4 + 64 = 84$

Sq. of Variations	Sum of Squares	Degrees of Freedom	Mean Sum of Sq.	$F = \frac{MSB}{MSW}$
B/W samples	$SSC = 160$	$V_1 = C-1 = 3-1=2$	$MSB = \frac{SSC}{V_1} = \frac{160}{2} = 80$	$F = \frac{80}{7}$
Within Sample	$SSE = 84$	$V_2 = n-C = (5-3=2)$	$MSW = \frac{SSE}{V_2} = \frac{84}{12} = 7$	$F = 11.42$

- Garden contains wala numerical
 → If the median of a distribution of given below is 28.5 then find n and y .

C.I	0-10	10-20	20-30	30-40	40-50	50-60
Frequency	5	x	20	15	10	5

$$\text{Total frequency} = 60$$

- probability marker (opt) in notes

Unit 3

formulas from ppt shared by many
 Numerical → Normal distribution from

^{opt} page no - 21

- sampling and its types. → 10 Marks.
- sampling - cluster VS stratified.
 (or area)

- prob sampling VS non-prob sampling
- procedure for hypothesis test
- z-test VS t-test.
- ANOVA → 10 marks

No. of accidents - chi-square - 10 mark. (goodness of fit)

- population VS mean

- Pmp of statistics → 2 marks

- what is p-value, critical region → 2m
 and similar definitions → 2marks

- teacher claim question.

~~Other~~

- In-person and online wala question

- 4000 hours wala - 10 marks

* * * * * - explain types of z-test - 1D, 2D, 1se
 - Poisson - 10 mark.