

Statistics-II

- Ankit Pangani

Unit 1 : Sampling distribution & estimation

Population : A set of similar items or events that we take for study of statistics. It is a collection of data from which a statistical sample is drawn for study.

Parameter : A parameter is a value (piece of data) that tells you something about whole population. It is more reliable than sample.

Statistic : A piece of data that describes only the sample of a population.

Sampling : The process of selecting a sample from a population for the study or research is called Sampling.

Sampling distribution : It is a probability distribution of statistic obtained from a larger number of samples drawn from specific population. It arrives through repeated sampling from a larger population.

Sampling distribution of mean :

The distribution of the values of the sample mean in repeated samples is called Sampling dist. of me (उद्दी पॉपुलेशन के अन्तर्गत संकेतिक समूहों का सामैपिंग मीन का प्रबलिटी डिस्ट्रिब्युशन)

$X \sim N(\mu, \sigma^2)$ from finite population
 * If sample is drawn from finite population
 (replacement is not done)
 $\sigma = \text{Population S.D}$ $N = \text{Popn size}$
 $\sigma^2 = \text{Population variance}$ $n = \text{Sample size}$

Date _____
 Page _____

$$\left[\text{mean } E(\bar{X}) = \mu \right], \left[V(\bar{X}) = \frac{\sigma^2}{n(N-n)} \right]$$

If sample is from infinite popn, $\frac{n}{N}$ can be neglected
 (replacement is possible)

$$\left[V(\bar{X}) = \frac{\sigma^2}{n} \right], \left[E(\bar{X}) = \mu \right], \left[S.E. = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}} \right]$$

Properties:

i) Sample mean (\bar{X}) is unbiased estimate of population mean (μ).

ii) The variance $V(\bar{X})$ depends on the sample size (n)
 & is equal to $\frac{\sigma^2}{n}$.

$$\text{Standard error} = \sqrt{V(\bar{X})} = \sqrt{\frac{\sigma^2}{n} \frac{(N-n)}{N-1}} = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

Sampling distribution of proportion

A proportion is the number of elements with a given characteristic divided by the total number of elements in that group.

$$\text{Population proportion } P = \frac{x}{N}, Q = 1 - P$$

$x = \text{no. of characteristics obtained in popn}$
 $N = \text{population size.}$

$$\text{sample proportion } P = \frac{x}{n}$$

$x = \dots \dots$ in sam
 $n = \text{sample size.}$

To study any qualitative study, we need to estimate population proportion. To estimate population proportion, we use sample proportion. The distribution followed by sample proportion is called sampling distribution of proportion.

* If population is infinite, (with replacement)

$$E(p) = p, \quad \text{Var}(p) = \frac{pq}{n}, \quad \text{Standard Error} [SE(p)] = \sqrt{\frac{pq}{n}}$$

* If population is finite, (without replacement)

$$E(p) = p, \quad \text{Var}(p) = \frac{pq}{n} \left(\frac{n-n}{n-1} \right)$$

$$\text{Standard Error } [SE(p)] = \sqrt{\frac{pq}{n} \left(\frac{n-n}{n-1} \right)}$$

Properties: same

$$\text{or } \sqrt{\frac{pq}{n} \left(\frac{n-n}{n-1} \right)}$$

Central limit theorem

It states that, as the sample size gets large enough, the sampling distribution of mean is approximately normally distributed. It is true regardless of the shape of distribution of the individual values in the population.

If $x_1, x_2, x_3, \dots, x_n$ is a random sample of size n of any population,

$$\text{sample mean } (\bar{x}) = \frac{x_1 + x_2 + \dots + x_n}{n}$$

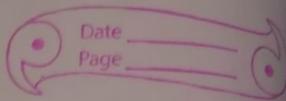
such that it is normally distributed, then

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

$$\text{Var}(\bar{x}) = \frac{\sigma^2}{n}$$

Q: 5

σ = population S.D (standard deviation)
 s = sample S.D



Important Formulas.

Statistic

Mean (σ known, pop'n size infinite)

$$S.E(\bar{x}) = \frac{\sigma}{\sqrt{n}} = \frac{\sigma}{\sqrt{n}}$$

Mean (σ known, pop'n size finite)

$$S.E(\bar{x}) = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

Mean (σ unknown, pop'n infinite)

$$S.E(\bar{x}) = \frac{s}{\sqrt{n}}$$

Mean (σ unknown, pop'n finite)

$$S.E(\bar{x}) = \frac{s}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$

Difference of means (σ 's are known)

$$S.E(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Difference of means (σ 's are unknown)

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

Proportion (pop'n size ~~finite~~ infinite)

$$S.E(P) = \sqrt{\frac{PQ}{n}}$$

Proportion (pop'n size ~~finite~~)

$$S.E(P) = \sqrt{\frac{PQ}{n} \left(\frac{n-n}{n-1} \right)}$$

Difference of props

$$S.E(P_1 - P_2) = \sqrt{\frac{P_1 Q_1}{n_1} + \frac{P_2 Q_2}{n_2}}$$

$$= \sqrt{PQ} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

where, N = population size.

n = sample size

σ = pop'n standard deviation

s = sample standard deviation

P = proportion of success, $Q = "1" - "failure"$

Inferential statistics

Descriptive statistics describes data (for example, a chart or graph) and inferential statistics allows you to make predictions from that data. With inferential statistics, you take data from samples & make generalization about a population.

estimation : The process of estimating the unknown population parameters like mean, variance, proportion from the sample statistics.

Estimator : The sample statistic which is used to estimate the unknown population parameters.
Ex: Sample mean

Estimate : The particular value taken by the estimator

→ Point estimation : A sample statistic (numerical value) is used to provide an estimate of population parameter

→ Interval estimation : Probable range is specified within which the value of parameters might be expected to lie.

→ good estimator criteria

1) Unbiasedness
2) Consistency

3) Efficiency
4) Sufficiency

* Confidence Interval (C.I) estimation of population mean (μ)

Large sample ($n > 30$)

$$C.I. = \bar{x} \pm z_{\alpha} S.E(\bar{x})$$

$$C.I. = \bar{x} \pm z_{\alpha} \frac{\sigma}{\sqrt{n}}$$

$$(I) \quad \bar{x} - z_{\alpha} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{\alpha} \frac{\sigma}{\sqrt{n}}$$

Small sample ($n < 30$)

$$q = \bar{x} \pm t_{\alpha, n-1} S.E(\bar{x})$$

where,
 z_{α} = signif. value of z
 n = sample size

\bar{x} = sample mean

σ = standard deviation

* Confidence Interval (C.I) estimation of population proportion (π)

Large sample ($n > 30$)

$(1-\alpha) \times 100\%$ confidence interval

$$C.I. = p \pm z_{\alpha} S.E(p)$$

$$C.I. = p \pm z_{\alpha} \sqrt{\frac{pq}{n}}$$

$$p \pm z_{\alpha} \sqrt{\frac{pq}{n} \left(\frac{N-n}{N-1} \right)}$$

If N is given

z_{α} = significance value of z
 where, p = sample proportion of success

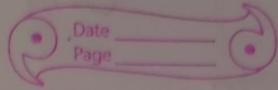
q = sample proportion of failure

n = sample size

N = population size

$$p = \frac{x}{n}, q = 1-p$$

(Ex characteristics taken from n sample)



- * Sample size determination for the estimate of popn mean (μ)

Minimum reqd sample size in estimating the population mean μ is

$$n = \frac{Z_{\alpha}^2 \sigma^2}{e^2}$$

where, Z_{α} = significance value of z
ex: 1.96 for 95% CI

σ = population s.d

e = margin of error

- * Sample size determination for the estimate of population proportion (π)

Minimum required sample size in estimating the population proportion is

$$n = \frac{Z_{\alpha}^2 \pi(1-\pi)}{e^2}$$

$$n = \frac{Z_{\alpha}^2 p \cdot q}{e^2}$$

where, n = sample size

Z_{α} = significance / critical value of z

π = population proportion, e = error margin

Note: * confidence level increases, sample size increases
 $C.I \propto n$

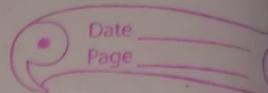
* Sampling error decreases, sample size increases
 $e \propto \frac{1}{\sqrt{n}}$

* sample size increases, width of CI decreases

It means for achieving higher no. of sample size, confidence level should increase & Sampling error should decrease

Q.S. 1. A ...

Unit-2: Testing of Hypothesis

Parametric

~~Hypothesis~~: A statistical hypothesis is a tentative statement or supposition about the estimated value of one or more population parameters.

Hypothesis: A hypothesis is a tentative theory or supposition provisionally adopted to explain certain facts & to guide in investigation.

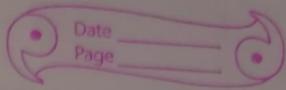
Non-parametric hypothesis: Statistical hypotheses about attributes is called non-parametric.

Note: If a hypothesis completely defines the population, it is a simple hypothesis, otherwise composite hypothesis.

Types of hypothesis

1) (Null hypothesis): The supposition about the population parameter is called null hypothesis. It is a hypothesis which is tested for possible rejection under assumption that it is true.

- No difference betn sample statistic & parameter denoted by H_0 . & set up as $H_0: \mu = \mu_0$



- 2) (Alternate hypothesis): Complementary to null hypothesis.
- difference between sample statistic & parameter.
 - denoted by H_1 & set up as $H_1: \bar{w} \neq w_0 \rightarrow$ two-tailed
- one tailed right $\bar{w} > w_0$] one tailed
one tailed left $\bar{w} < w_0$]

Errors in hypothesis testing:

On the basis of accepting or rejecting the hypothesis two types of errors are produced.

- 1) Type I error (null hyp. true $\bar{g} \& H_0$ not rejected)
- It is the error of rejecting null hypothesis when it is true. Its probability is α (Level of significance)

- 2) Type II error (null hyp. false $\bar{g} \& H_0$ accepted)
- It is the error of accepting null hypothesis when it is false (ie. when alternative hypothesis is true). Its probability is β . [$\beta = \text{Prob}(\text{accepting } H_0 | H_1 \text{ is true})$]

Power of test ($1-\beta$) (first explain about β)

The probability of rejecting H_0 when H_1 is false is probability of correct decision which is called ~~prob~~ power of test.

i.e. It is probability of not making type II error / $\text{Prob. of accepting } H_1 \text{ when } H_1 \text{ is true}$.

Test statistic

It is the statistic based upon approximate probability distribution. It is used to decide whether to accept or reject null hypothesis. Commonly used test statistics:

1) Z-test: we use z-distribution under the normal curve for large sample ($n > 30$)

The z-test statistic is $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$ ~ $N(0, 1)$

2) t-test: we use t-distribution for small sample ($n < 30$).

The t-test statistic is $t = \frac{\bar{x} - \mu}{S.E(\bar{x})}$ ~ t dist. with $n-1$ degrees of freedom

Level of significance (L.O.S.)

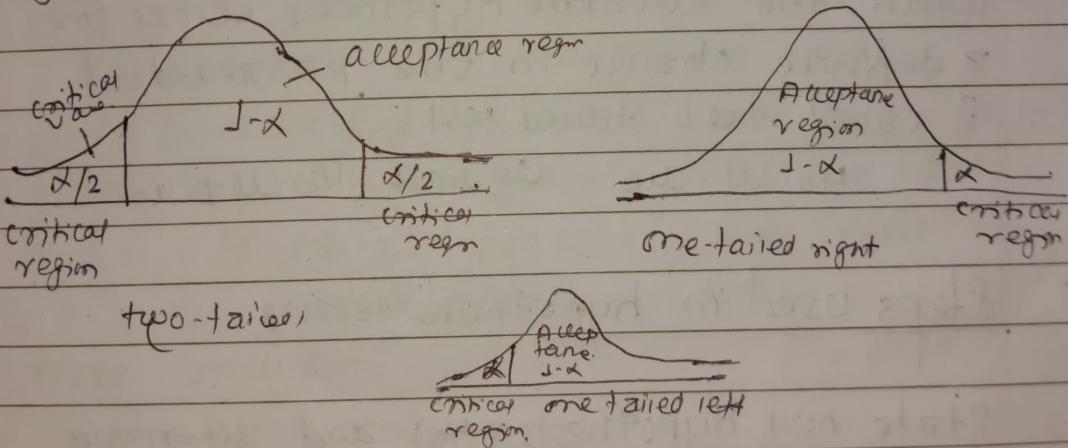
The two errors (type-I & type-II) are inversely related. Both errors cannot be minimized at the same time. Usually we fix type I & minimize type II). The maximum size of type I error prepared in testing of hypothesis is called L.O.S. Denoted by α .

→ commonly used L.O.S are 5%, 1%. The L.O.S should be chosen on basis of power of test. If P.O.T is too low then high L.O.S should be chosen, i.e. 10%, 20%.

Critical region.

Also called the rejection region. The set of all possible values of statistic is divided into two regions, one leading to rejection & others leading to acceptance of null hypothesis (H_0). The division is based on L.O.S and alternate hypo (H_1).

- region which leads to rejection of H_0 = rejection region ($\cup \cap$)
- region which " " " acceptance " H_0 = acceptance region ($\cap \cup$)
- If test statistic falls in rejection region, then H_0 is rejected, otherwise accepted.



Critical value

The value of statistic which separates critical region and acceptance region is called critical value.

Degree of freedom (No. of independent choice / choice selection)

The no. of independent variates which make up statistic is called d.f. If sample size = n & restriction = 1

Then, $d.f = n - 1$

* One-tailed test (for directional hypothesis)
A test of statistical hypothesis in which the alternative hypothesis H_1 looks for a definite increase (right tail) or definite decrease (left tail) in parameter is called

<u>NULL</u>	<u>Alternate</u>	
$H_0: \mu = \mu_0$ against	$H_1: \mu > \mu_0$	(right tail)
$H_0: \mu = \mu_0$ against	$H_1: \mu < \mu_0$	(left tail)

* Two-tailed test (for non-directional hypothesis)
A test of statistical hypothesis in which the alternate hypothesis looks for a definite change in the parameter is called two tailed test.

$$H_0: \mu = \mu_0 \text{ against } H_1: \mu \neq \mu_0$$

* Steps used in hypothesis testing

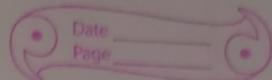
Step-1: State null hypothesis (H_0) and alternate hypothesis (H_1)

Step-2: Select the appropriate test statistics.

Step-3: State level of significance. and calculate the required value for the test. under H_0 .

Step-4: Find the critical value (tabulated value).

Step-5: Make decision by comparing the calculated value & tabulated value.



reject.

Step-5: Make conclusion: i.e. either accept H_0 or ~~reject H_0~~ .

* Procedure of testing significance of mean for large sample ($n \geq 30$)

Step-1: State H_0 & H_1 ($H_0: \mu = \bar{x}_y_2 / \bar{m} = u_0$, $H_1: \mu > \bar{x}_y_2$)

Step-2: Select Z-test & find test statistic under H_0

$$Z = \frac{\bar{x} - u}{\frac{\sigma}{\sqrt{n}}}, \quad \bar{x} = \text{sample mean}$$

u = population/hypothesised mean

σ = population standard deviation

find $|Z|$

n = sample size

Step-3: State L.O.S & find tabulated value of Z tab $|z|$ (critical value)

Step-4: Make decision by comparing $|Z|$ & $tab|z|$

If $|Z| < tab|z|$ then accept H_0

reject otherwise.

Step-5: Make conclusion.

* Procedure of testing significance of mean for small sample ($n < 30$)

$$\left\{ \begin{array}{l} t = \frac{\bar{x} - u}{\frac{s}{\sqrt{n}}} \text{ use this when sample SD is already given} \\ \text{or } t = \frac{\bar{x} - u}{\frac{s}{\sqrt{n-1}}} \text{ otherwise use } s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} \end{array} \right.$$

Step-1: State H_0 & H_1 .

Step-2: Select t-test and find test statistic under H_0

$$t = \frac{\bar{x} - u}{\frac{s}{\sqrt{n-1}}} \quad \text{or} \quad t = \frac{\bar{x} - u}{\frac{s}{\sqrt{n-1}}} \quad \text{or} \quad t = \frac{\bar{x} - u}{\frac{s}{\sqrt{n-1}}} \quad \left(\frac{s}{\sqrt{n-1}} \right) \quad \text{if } s \text{ is given}$$

where, \bar{x} = Sample mean
 μ = Population mean
 $s/\delta s$ = sample S.D. = $\sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$

σ = population S.D.

Step-3: State L.O.S & find tabulated value of t
 Step-4: Find critical value, $|t|_{\text{tab}}$ and compare with $|t|_{\text{cal}}$

Step-5: Make decision: If $|t|_{\text{cal}} < |t|_{\text{tab}}$ accept H_0
 reject otherwise.

Step-6: Make conclusion.

* procedure for testing significance of difference between two means for small sample

Step-1: State H_0 and H_1 ($H_0: \mu_1 = \mu_2$) $H_1: \mu_1 \neq \mu_2, \mu_1 < \mu_2, \mu_1 > \mu_2$

Step-2: Select t-test and find test statistic under H_0

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

$$\text{or } t = \frac{\bar{x}_1 - \bar{x}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

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

$$S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

where,

S_p^2 = pooled / combined variance

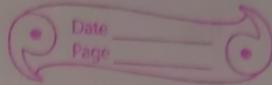
s_1 = 1st sample SD

n_1 = 1st sample size

s_2 = 2nd sample SD

n_2 = 2nd sample size

Note: confidence level = $1 - \alpha$
 $95\% = 1 - \alpha$
 $0.95 = 1 - \alpha$
 $\alpha = 0.05$



Here,

$$Sp^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2 - 2} \quad \text{or} \quad Sp^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

use If s_1 & s_2 are not given.

Step 1: State Level of significance (α)

Step 2: Find critical value / table value of t

Step 3: Make decision by comparing $|cal|t|$ and $|table|t|$
 If $|cal|t| < |table|t|$, accept H_0 , reject otherwise.

Step 4: Make Conclusion

* procedure for testing significance of proportion for large sample ($n > 30$)

Step 1: $H_0: P = Xyz$, $H_1: P \neq, P >, P < Xyz$

Step 2: Select Z-test & find test statistics under H_0 .

$$Z = \frac{\bar{P}_S - P}{\sqrt{\frac{P_S(1-P_S)}{n}}} \quad \bar{P}_S = \text{sample proportion} = \frac{x}{n} \quad P = \text{population proportion}$$

$$Q_S = 1 - P_S$$

$$n = \text{sample size}$$

Step 3: $" "$

Step 4: $" "$

Step 5: If $|cal|z| < |table|z|$, accept H_0 or reject H_0

Step 6: $" "$

$$= 0.4$$

* Procedure of testing significance of difference between two proportions for large samples ($n > 30$)

Steps: State H_0 and H_1 ($H_0: P_1 = P_2$, $H_1: P_1 \neq P_2$ or $P_1 > P_2$, $P_1 < P_2$)

Step 2: Select z-test x "

$$z = \frac{p_1 - p_2}{\sqrt{pq(\frac{1}{n_1} + \frac{1}{n_2})}}$$

p_1 & p_2 are sample prop
 $p_1 = \frac{x_1}{n_1}$, $p_2 = \frac{x_2}{n_2}$

or

$$P = \frac{x_1 + x_2}{n_1 + n_2}, Q = 1 - P$$

$$z = \frac{p_1 - p_2}{\sqrt{\frac{PQ}{n_1} + \frac{PQ}{n_2}}}$$

P, Q are popⁿ prop's

when p_1 and p_2 are given.

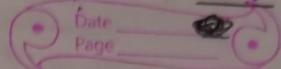
Step-3: " " " "

Step-4: " "

Step-5: Make conclusion $\Rightarrow |z| \geq |z|_{\text{table}}$
accept H_0 else reject H_0 :

95% CI, ~~for~~: For z-test, if $\alpha = 0.05$, for two-tailed $Z_{\alpha} = \frac{1-\alpha}{2}$

95% CI, for t-test, if $\alpha = 0.05$, for two-tailed, $t_{\alpha} = \frac{1-\alpha}{2}$



Note: for z-table (not standard normal dist.)
if $\alpha = 0.05$

for one-tailed \Rightarrow see $(1-\alpha)$ for Z_{α} (critical value)

for two-tailed \Rightarrow see $(1-\alpha/2)$ for $Z_{\alpha/2}$ (critical value)
if $|Z_{\text{cal}}| > |Z_{\text{table}}|$ reject H_0 , otherwise accept H_0

for t-test (table in LAPP)

for one-tailed \Rightarrow see α and $(n-1)$ df

for two-tailed \Rightarrow see α and $(n-1)$ df.

If $|t_{\text{cal}}| > t_{\text{table}}$, then reject H_0 , otherwise accept H_0

for p-value test

for one-tailed \Rightarrow P-value = $\text{Prob}\{Z > |Z_{\text{cal}}|\}$

for two-tailed \Rightarrow P-value = $2\text{Prob}\{Z > |Z_{\text{cal}}|\}$

if P-value $< \alpha \Rightarrow$ reject H_0
otherwise accept H_0

for paired t-test

X = before, Y = after

$H_0: \mu_x = \mu_y$

$H_1: \mu_x \neq \mu_y, \mu_x < \mu_y, \mu_x > \mu_y$.

Test statistic: $t = \frac{\bar{d}}{\frac{s_d}{\sqrt{n}}}$ where $\bar{d} = \frac{\sum d}{n}$.
 $s_d = \sqrt{\frac{1}{n-1} \sum (d - \bar{d})^2}$

standard dev = $s_d = \sqrt{\frac{1}{n-1} \sum (d - \bar{d})^2}$

or $s_d = \sqrt{\frac{\sum (d - \bar{d})^2}{n-1}}$

$$= \sqrt{\frac{1}{n-1} \left\{ \sum d^2 - \frac{(\sum d)^2}{n} \right\}}$$