

DIGITAL NOTES ON PROBABILITY & STATISTICS

**B.TECH II YEAR - I SEM
(2018-19)**



DEPARTMENT OF INFORMATION TECHNOLOGY

**MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY
(Autonomous Institution – UGC, Govt. of India)**

(Affiliated to JNTUH, Hyderabad, Approved by AICTE - Accredited by NBA & NAAC – ‘A’ Grade - ISO 9001:2015 Certified)
Maisammaguda, Dhulapally (Post Via. Hakimpet), Secunderabad – 500100, Telangana State, INDIA.



MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY
DEPARTMENT OF INFORMATION TECHNOLOGY

II B.Tech I Sem

L	T/P/D	C
4	1 / - / -	3

(R17A0024)PROBABILITY AND STATISTICS

Objectives: To learn

- Understand a random variable that describes randomness or an uncertainty in certain realistic situation. It can be either discrete or continuous type.
- In the discrete case, study of the binomial and the Poisson random variables and the normal random variable for the continuous case predominantly describe important probability distributions. Important statistical properties for these random variables provide very good insight and are essential for industrial applications.
- Most of the random situations are described as functions of many single random variables. The objective is to learn functions of many random variables, through joint distributions.
- The types of sampling, Sampling distribution of means, Sampling distribution of variance, Estimations of statistical parameters, Testing of hypothesis of few unknown statistical parameters.
- The mechanism of queuing system, The characteristics of queue, The mean arrival and service rates, The expected queue length, The waiting line, The random processes, The classification of random processes, Markov chain, Classification of states, Stochastic matrix (transition probability matrix), Limiting probabilities, Applications of Markov chains.

UNIT -1 : Random variable and Probability distributions

Random Variables

Single and multiple Random variables -Discrete and Continuous. Probability distribution function, mass function and density function of probability distribution. mathematical expectation and variance.

Probability distributions: Binomial distribution – properties, mean and variance, Poisson distribution – properties, mean and variance and Normal distribution – properties, mean and variance

UNIT -2 :Correlation and Regression

Correlation -Coefficient of correlation , Rank correlation, Regression- Regression Coefficients , Lines of Regression.

UNIT -3 : Sampling Distributions and Statistical Inferences

Sampling: Definitions of population ,sampling ,statistic ,parameter-Types of sampling – Expected values of sample mean and variance,Standard error- Sampling distribution of means and variance

Parameter Estimations : likelihood estimate , interval estimate.

Testing of hypothesis: Null and Alternative hypothesis-Type I and Type II errors, Critical region – confidence interval – Level of significance,One tailed and Two tailed test

Large sample Tests: i) Test of significance of single mean and equality of means of two samples(cases of known and unknown variance whether equal or unequal)

ii) Tests of significance difference between sample proportion and population proportion and difference between two sample proportions

UNIT -4 : Exact Sampling Distributions(Small samples)

Exact Sampling Distributions(Small samples) Student t- distribution - properties

i)Test of significant difference between sample and population mean

ii)Test of difference between means of two small samples(independent and dependent samples)

F- distribution - properties –test of equality of two population variances

Chi-square distribution -properties –i)Test of goodness of fit

ii)Test of independence of attributes

UNIT-5 : Queuing Theory and Stochastic process

Queuing Theory

Structure of a queuing system its characteristics-Arrival and service process-Pure Birth and Death process Terminology of queuing system -Queuing model and its types-M/M/1 model of infinite queue (without proofs)and M/M/1 model of finite queue (without proofs).

Stochastic Process

Introduction to stochastic process-classification and methods of description of Random process i.e,stationary and non-stationary Average values of single and two or more random process

Markov process, Markov chain, Examples of Markov chains, Stochastic matrix.

TEXT BOOKS:

1. Probability and Statistics by T.K..V Iyengar& B.Krishna Gandhi S.Ranganatham,MVSSAN Prasad. S Chand Publishers.
2. Fundamentals of Mathematical Statistics by SC Gupta and V.K. Kapoor

REFERENCES :

- 1.Higher Engineering Mathematics By Dr.B.S.Grewal, Khanna Publishers
2. Probability and Statistics for Engineers and Scientists by Sheldon M.Ross,Academic Press.

Outcomes:

- Students would be able to identify distribution in certain realistic situation. It is mainly useful for circuit as well as non circuit branches of engineering. Also able to differentiate among many random variables involved in the probability models. It is quite useful for all branches of engineering.
- The student would be able to calculate mean and proportions(small and large samples)and to make important decisions from few samples which are taken out of unmanageably huge populations.It is mainly useful for non-branches of engineering.

- The student would be able to find the expected queue length, the ideal time the traffic intensity and the waiting time. these are very useful tools in many engineering and data management problems in the industry. it is useful for all branches of engineering.
- The student would be able to understand about the random process, markov process and markov chains which are essentially models of many time dependent processes such as signals in communications, time series analysis, queuing systems. The student would be able to find the limiting probabilities and the probabilities in n^{th} state. It is quite useful for all branches of engineering.



MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY
DEPARTMENT OF INFORMATION TECHNOLOGY

INDEX

S. No	Unit	Topic	Page no
1	I	Introduction to Random variables	6
2	I	Probability Distribution	10
3	I	Multiple Random Variables	16
4	I	Problems of Unit 1	21
5	II	Correlation	40
6	II	Regression	58
7	III	Sampling Distribution	72
8	III	Estimation	80
9	III	Testing of Hypothesis	81
10	IV	Small Samples	106
11	IV	t- distribution	115
12	IV	F- distribution & Chi-Square distribution	121
13	IV	Related Problems	126
14	V	Queuing Theory	140
15	V	Stochastic Process	154

NOTE:-List only main topics

Probability function of a Discrete random variable :

If for a discrete r.v 'X', the real valued function $P(x)$ is such that $P(X=x) = p(x)$ then $p(x)$ is called Probability function or probability mass function of a discrete r.v 'X'.

Probability Distribution function .

P.D.F associated with X is the probability that the outcome of an experiment will be one of the outcomes for which $X(s) \leq x, x \in R$

i.e $F_x(x) = P(X \leq x) = P\{S : X(S) \leq x\}, -\infty < x < \infty$
is called Distribution function .

Properties of Distribution function

1. If F distribution function of X , and if $a < b \Rightarrow P(a < X \leq b) = F(b) - F(a)$

2. $P(a \leq X \leq b) = P(X=a) + [F(b) - F(a)]$

3. $P(a < X < b) = [F(b) - F(a)] - P(X=b)$

4. $P(a \leq X < b) = [F(b) - F(a)] - P(X=b) + P(X=a)$

Note: $0 \leq F(x) \leq 1$.

Discrete Probability Distribution (Probability mass function)

It is the set of its possible values together with their respective probabilities.

Let X be a discrete r.v with possible outcomes x_1, x_2, x_3, \dots

then their probabilities $p_i = P(X=x_i) = P(x_i)$ for $i=1, 2, 3, \dots$

if $p(x_i) > 0$ & $\sum_{i=1}^n p(x_i) = 1$ then the function ' $p(x)$ ' is called the

'probability mass function' of r.v X & $\{(x_i, p(x_i))\}$ for $i=1, 2, \dots$ is called Discrete Probability Distribution.

Eg: Tossing a Coin two times with random variable $X(s) = \text{no. of heads} \in \{0, 1, 2\}$
then $P(X=0) = \frac{1}{4}$; $P(X=1) = \frac{1}{2}$, $P(X=2) = \frac{1}{4}$

x_i	0	1	2
$P(x_i)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

$\Rightarrow \sum p(x_i) = 1 \& p(x_i) > 0$.

thus total probability '1' is distributed into 3 parts as $\frac{1}{4}, \frac{1}{2}$ & $\frac{1}{4}$.

Note: 1) Frequency distribution tells how total frequency distributed among different values of the variable .

2) Probability distribution tells how total probability '1' is distributed among the values which the random variable can take .

2

Cumulative Distribution Function of a Discrete R.V.

$F(x) = P(X \leq x) = \sum_{i=1}^{\infty} P(x_i)$ where x is any integer.
 Then $F(x)$ is called cumulative distribution function of X .
 i.e. $F(x) = \begin{cases} 0, & -\infty < x < x_1 \\ P(x_1), & x_1 \leq x < x_2 \\ P(x_2), & x_2 \leq x < x_3 \\ \vdots & \vdots \\ P(x_1) + P(x_2) + \dots + P(x_n), & x_n \leq x < \infty \end{cases}$

Expectation of Discrete Probability Distribution.

Let X be a r.v. with x_1, x_2, \dots, x_n with probabilities p_1, p_2, \dots, p_n .

then

Mathematical Expectation or Expected value of X is defined as "sum of products of different values of X & their corresponding probabilities".

$$\text{i.e. } E(X) = \sum_{i=1}^n x_i p_i$$

$$\text{In General } E(g(x)) = \sum_{i=1}^n p_i g(x_i)$$

$$\text{Results: 1. } E(X) = \mu \quad \left\{ \because \text{Mean} = \frac{\sum p_i x_i}{\sum p_i} = \sum p_i x_i \text{ as } \sum p_i = 1 \right\}$$

$$2. \quad E(X+k) = E(X)+k \quad (k \text{ is a constant})$$

$$3. \quad E(kX) = k E(X)$$

$$4. \quad E(aX+b) = aE(X)+b$$

$$5. \quad E(X+Y) = E(X)+E(Y) \text{ provided } E(X), E(Y) \text{ exists.}$$

$$6. \quad E(X-\bar{X}) = 0$$

$$7. \quad E(XY) = E(X) \cdot E(Y) \text{ if } X, Y \text{ are two independent r.v.s.}$$

$$8. \quad E(\frac{1}{X}) \neq \frac{1}{E(X)}$$

Mean: The mean ' μ ' of the distribution function is given by

$$\mu = \frac{\sum p_i x_i}{\sum p_i} = \sum p_i x_i$$

Variance: Variance of the probability distribution of r.v. X is defined as

$$\text{Var}(X) = E[(X - E(X))^2]$$

$$\text{Since variance of } X = V(X) = \sum_{i=1}^n (x_i - \mu)^2 p_i = \sum_{i=1}^n p_i x_i^2 - \mu^2 \\ = E(X^2) - [E(X)]^2$$

$$\text{Also } V(X) = E(X^2) - [E(X)]^2$$

Standard Deviation $S.D = \sigma = \sqrt{\text{Var}(X)} = \sqrt{E(X) - [E(X)]^2}$

Results: 1. $\text{Var}(k) = 0$ where k is constant $\left\{ \begin{array}{l} \text{mean = constant} \\ \Rightarrow X - \mu = 0 \end{array} \right.$

$$2. \text{Var}(kx) = k^2 \text{Var}(x)$$

$$3. \sqrt{X+k} = \sqrt{X}$$

* 4. $\sqrt{ax+b} \neq a^2 \sqrt{X}$ where a, b are constants

Let $Y = ax+b$

$$\therefore E(Y) = E(ax+b)$$

$$= aE(X) + b$$

$$\text{Consider } Y - E(Y) = ax+b - aE(X) - b = a[X - E(X)]$$

On taking expectation, we get

$$E(Y - E(Y))^2 = a^2 E(X - E(X))^2$$

$$\Rightarrow \text{Var}(Y) = a^2 \text{Var}(X)$$

$$\Rightarrow \sqrt{a^2 \text{Var}(X)} = a \sqrt{\text{Var}(X)}$$

Hence Proved.

Do Problems.

Continuous Probability Distribution: (Probability Density function)

Defn: By considering small interval $[x - \frac{dx}{2}, x + \frac{dx}{2}]$ of length dx around the point x . Let $f(x)$ be any continuous function of x so that $f(x)dx$ represents the probability that the variable X falls in that interval. i.e. $P(x - \frac{dx}{2} \leq X \leq x + \frac{dx}{2}) = f(x)dx$

then $f(x)$ is called the probability density function & the curve $y=f(x)$ is known as probability density curve.

The probability for a variate value to fall in the finite interval (a, b) is $\int_a^b f(x)dx$ which represents the area below the curve $y=f(x)$, $x=a$ & $x=b$.

Properties of density function

$$(i) f(x) \geq 0 \quad \forall x \in \mathbb{R}$$

$$(ii) \int_{-\infty}^{\infty} f(x)dx = 1$$

(iii) $P(E) = \int_E f(x)dx$ is well defined for any event E .

(iv) Since continuous variable associate with intervals, the probability at a particular point is always zero. Thus

$$P(a < X \leq b) = P(a \leq X < b) = P(a < X < b) = P(a \leq X \leq b) = F(b) - F(a)$$

* Note: Since continuous variable associate with intervals, the probability at a particular point is always zero. Thus

Binomial Distribution

It was discovered by James Bernoulli in the year 1700 & is a discrete probability distribution.

The conditions for the applicability of a B.D are :

- (i) There are 'n' independent trials (ie-trials are repeated under identical conditions)
- (ii) There are only 2 possible outcomes for each trial. Success 'p' & failure 'q' .
- (iii) The trials are independent, ie the probability of an event in any trial is not affected by the results of any other trial .
- (iv) The probability of success in each trial remains constant & does not change from trial to trial .

Defⁿ: A rv.X has a B.D if it assumes only non-negative values & its probability density function is given by

$$P(X=r) = p(r) = \begin{cases} {}^n C_r p^r q^{n-r} & ; r=0,1,2,\dots,n ; q=1-p \\ 0 & ; \text{otherwise} \end{cases}$$
$$= b(r; n, p)$$

Here n, p are called parameters of the distribution as they are the 2 independent constants. 'n' is sometime known as "degree of the distribution".

- e.g.s of B.D : (i) No. of defective bolts in a box containing 'n' bolts .
(ii) No. of post graduates in a group of 'n' men , etc .

The Binomial Distribution function is:

$$F_x(x) = P(X \leq x) = \sum_{r=0}^n {}^n C_r p^r q^{n-r} .$$

Binomial Frequency Distribution: The possible no. of success and their frequencies is called a Binomial frequency Distribution .

If 'n' independent trials constitute one experiment and this exp. is repeated 'N' times, then the frequency of 'r' successes is
 $N \cdot {}^n C_r p^r q^{n-r}$.

\therefore In 'N' sets of 'n' trials the theoretical frequencies of 0, 1, 2, ..., r, ..., n successes are given by the terms of expansion of $N(q+p)^n$.

Note ① The probabilities of 0, 1, 2, ..., r, ..., n successes in 'n' trials are given by the terms of the binomial expansion of $(q+p)^n$.

$$\text{ie } (q+p)^n = q^n + {}^n C_1 q^{n-1} p + {}^n C_2 q^{n-2} p^2 + \dots + {}^n C_r q^{n-r} p^r + \dots + p^n.$$

Here the probability of exactly 'r' successes is ${}^n C_r p^r q^{n-r}$.

② The probability of no success in 'n' trials is $= q^n$

" " " all successes " " " = p^n .

$$\begin{aligned} " " " \text{ " at least One success" " " " } &= {}^n C_1 q^{n-1} p + {}^n C_2 q^{n-2} p^2 \\ &+ \dots + p^n = \underline{\underline{1 - q^n}}. \end{aligned}$$

~~Do problems.~~

Constants of B.D :

$$\begin{aligned} \text{(1) Mean of } X : \mu &= E(X) = \sum_{r=0}^n r P(r) = \sum_{r=0}^n r \cdot {}^n C_r p^r q^{n-r} \\ \Rightarrow \mu &= 1 \cdot {}^n C_1 p q^{n-1} + 2 {}^n C_2 p^2 q^{n-2} + \dots + n {}^n C_n q^{n-n} p^n \\ &= n p q^{n-1} + \frac{2 n (n-1)}{2!} p^2 q^{n-2} + \dots + n p^n \\ &= n p \left[q^{n-1} + \frac{n(n-1)}{2!} p q^{n-2} + \frac{n(n-1)(n-2)}{3!} p^2 q^{n-3} + \dots + p^n \right] \\ &= n p \left[q^{n-1} + \frac{(n-1)p q^{n-2}}{2!} + \frac{(n-1)(n-2)p^2 q^{n-3}}{3!} + \dots + p^n \right] \\ &= n p \left[{}^n C_0 q^{n-1} p + {}^n C_1 p q^{n-2} + {}^n C_2 p^2 q^{n-3} + \dots + {}^n C_{n-1} p^{n-1} \right] \\ &= n p [q + p]^{n-1} \quad . \quad [\because \mu = np] \end{aligned}$$

2. Variance of B.D:

$$V(X) = E(X^2) - [E(X)]^2 = npq$$

$$\begin{aligned}
 \text{Proof: } V(X) &= E(X^2) - [E(X)]^2 \\
 &= \sum_{r=0}^n r^2 p(r) - \left[\sum_{r=0}^n r p(r) \right]^2 \\
 &= \sum_{r=0}^n [r(r-1) + r] p(r) - \mu^2 \\
 &= \sum_{r=0}^n r(r-1)p(r) + \sum_{r=0}^n r p(r) - \mu^2 \\
 &= \sum_{r=0}^n r(r-1)p(r) + \mu - \mu^2 \\
 &= 2(2-1)^n C_2 p^2 q^{n-2} + 3(3-1)^n C_3 p^3 q^{n-3} + \dots + \\
 &\quad n(n-1)^n C_n p^n q^{n-n} + \mu - \mu^2 \\
 &= 2^n C_2 p^2 q^{n-2} + 6^n C_3 p^3 q^{n-3} + \dots + n(n-1) p^n + \mu - \mu^2 \\
 &= \frac{n(n-1)}{2} p^2 q^{n-2} + \frac{n(n-1)(n-2)}{3 \times 2} p^3 q^{n-3} + \dots + n(n-1) p^n \\
 &= n(n-1) p^2 [q^{n-2} + (n-2)pq^{n-3} + \dots + p^{n-2}] \\
 &\quad + \mu - \mu^2 \\
 &= n(n-1) p^2 \left[q^{n-2} + C_1 p q^{n-3} + \dots + C_{n-2} p^{n-2} \right] \\
 &= n(n-1) p^2 [q + p]^{n-2} + \mu - \mu^2 \\
 &= n(n-1) p^2 (\because q + p = 1) + \mu - \mu^2 \\
 &= n^2 p^2 - np^2 + np - p^2 \quad \{\because \mu = np\} \\
 &= np(1-p) \\
 &= npq \quad \{\because p+q=1 \Rightarrow q=1-p\}
 \end{aligned}$$

$$\therefore \text{Var}(X) = \sigma^2 = npq$$

Poisson Distribution

7

S. D. Poisson (Simeon Denis Poisson) (1837) introduced Poisson Distribution as a rare distribution of rare events. i.e. the events whose probability of occurrence is very small but the no. of trials which could lead to the occurrence of the event, are very large.

Conditions of P.D : (i) The no. of trials 'n' is large.

(ii) The probability of success 'p' is very small.

(iii) $np = \lambda$ is finite.

Examples of P.D :

(1) The no. of printing mistakes per page in a large text.

(2) The no. of cars passing a certain point in 1 minute.

Definition: A r.v X is said to follow P.D if it assumes only non-negative values & its probability density function is given by $p(x, \lambda) = P(X=x) = \begin{cases} e^{-\lambda} \frac{\lambda^x}{x!}, & x=0,1,2,\dots \\ 0, & \text{otherwise} \end{cases}$

Here $\lambda > 0$ is called parameter of the distribution.

Constants of Poisson Distribution :

$$(1) \text{ Mean } \mu = E(X) = \lambda \quad \because \sum_{x=0}^{\infty} x p(x) = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{x^x}{x!}$$

$$\text{Proof : } E(X) = \sum_{x=0}^{\infty} x p(x) = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{x^x}{x!} \\ = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!(x-1)!} = e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} \\ = e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda \cdot \lambda^{x-1}}{(x-1)!} = \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!}$$

$$\text{Take } x-1=y \quad = \lambda e^{-\lambda} \left[1 + \lambda + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \dots \right] \\ = \lambda e^{-\lambda} [e^\lambda] \\ = \lambda$$

$$\therefore \mu = \lambda$$

(ii) Variance of P.D.

$$V(X) = \lambda$$

$$\text{Proof: } \text{Var}(X) = E(X^2) - [E(X)]^2 = E(X^2) - \mu^2 = E(X^2) - \lambda^2.$$

$$E(X^2) = \sum_{x=0}^{\infty} x^2 p(x) = \sum_{x=0}^{\infty} x^2 \frac{e^{-\lambda} \lambda^x}{x!} = \sum_{x=1}^{\infty} x \frac{e^{-\lambda} \lambda^x}{(x-1)!}$$

$$= e^{-\lambda} \sum x \frac{\lambda^x}{(x-1)!} = e^{-\lambda} \sum [(x-1)+1] \frac{\lambda^x}{(x-1)!}$$

$$= e^{-\lambda} \sum (x-1) \frac{\lambda^x}{(x-1)!} + e^{-\lambda} \sum \frac{\lambda^x}{(x-1)!}$$

$$= e^{-\lambda} \sum_{x=2}^{\infty} \frac{\lambda^{x-2} \lambda^2}{(x-2)!} + \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!}$$

$$= \lambda^2 e^{-\lambda} \sum_{x=2}^{\infty} \frac{\lambda^{x-2}}{(x-2)!} + \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!}$$

$$= \lambda^2 e^{-\lambda} (\lambda^2) + \lambda e^{-\lambda} (\lambda)$$

$$= \lambda^2 e^{-\lambda} + \lambda e^{-\lambda} = \lambda + \lambda \Rightarrow E(X^2) = \lambda^2 + \lambda$$

$$\therefore \text{Var}(X) = E(X^2) - \lambda^2 = \lambda^2 + \lambda - \lambda^2 = \lambda \Rightarrow \boxed{\text{Var}(X) = \lambda}$$

Note. S.D. $\sigma = \sqrt{\lambda}$

Recurrence Relation:

$$P(x+1) = \frac{\lambda}{x+1} P(x)$$

Note: When n is large say greater than 30 & p is very small say less than 0.1 then B.D can be approximated by P.D.

Normal Distribution :

It is a continuous distribution. N.D was first discovered by English Mathematician De-Moivre (1667-1745) in 1733 & further refined by French Mathematician Laplace (1749-1827) in 1774 & independently by Karl Friedrich Gauss (1777-1851).

It is also known as Gaussian Distribution.

It is another limiting form of Binomial Distribution. It is found so often in real life that it is called Normal Distribution, the name which is commonly used.

Defn A r.v X is said to have a N.D if its p.d.f is given by $f(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$, $-\infty < x < \infty$, $\sigma > 0$,

μ is Mean & σ is Standard deviation of x . $-\infty < \mu < \infty$.

μ, σ are called parameters of N.D.

A linear combination of independent normal variables is also a normal variate.

Chief Characteristics of N.D : The graph of Normal Distribution in xy -plane is known as normal curve.

1. The graph of Normal Distribution in xy -plane is known as normal curve.
2. The curve is a bell-shaped curve & symmetrical w.r.t mean ' μ ', ie The two tails on right & left sides of the mean ' μ ' extends to infinity.
3. Area under the normal curve represents total population.
4. Normal curve is Unimodal (ie has only one max. pt.) At $x = \mu$: mean = Median = Mode (as distribution is symmetrical).
5. x -axis is an asymptote to the curve.
6. The points of inflexion of the curve are at $x = \mu \pm \sigma$ & curve changes from concave to convex at $x = \mu \pm \sigma$, to $x = \mu - \sigma$.

7. Area under the normal curve is distributed as :

- (i) Area of normal curve b/w $\mu-\sigma$ & $\mu+\sigma$ is 68.27%.
- (ii) " " " " " " " $\mu-2\sigma$ & $\mu+2\sigma$ is 95.43%.
- (iii) " " " " " " " $\mu-3\sigma$ & $\mu+3\sigma$ is 99.73%.

Note : The total area bounded by the curve & x -axis is One.

$$\text{i.e. } \int_{-\infty}^{\infty} f(x) dx = 1.$$

$P(a < x < b)$ = Area under normal curve b/w the vertical lines $x=a$ & $x=b$.

$$= \int_a^b f(x) dx.$$

Multiple Random Variables.

In many practical problems several r.v.s interact with each other. Multiple R.Vs help us to determine the joint statistical properties like mean, variance etc.

We study about 2-dimensional r.v.s, which can be easily extended to multiple r.v.s.

A 2-dimensional r.v. is denoted by (X, Y) where the random vector (X, Y) is outcome of a trial which occurs in pairs i.e. $X=x, Y=y$.

Defn: An 'n' dimensional random vector (vector of r.v.s) is a function from sample space S into R^n (n -dimensional Euclidean Space).

Multiple r.v.s are also of 2 types: ① Discrete & ② Continuous.

It is called Discrete if (X, Y) assume only finite no. of pairs.

Discrete Multiple Random Variables:

Joint Probability Mass function (Joint Probability function)

If (X, Y) is a discrete 2-dimensional r.v. then the function $f(x, y)$ from R^2 into R is called JPMF & is given by
$$f(x, y) = P(X=x, Y=y) = \sum_{i} \sum_{j} P(x_i, y_j)$$

Note: ① $0 \leq P(x_i, y_j) \leq 1$. ② $\sum_{i} \sum_{j} P(x_i, y_j) = 1$.

Properties: ① $P(X=x_i) = p(x_i) = \sum_j P(x_i, y_j)$

② $P(Y=y_j) = p(y_j) = \sum_i P(x_i, y_j)$

③ $P(x_i) \geq P(x_i, y_j)$ for any j

④ $P(y_j) \geq P(x_i, y_j)$ for any i

Joint Cumulative Distribution Function .

Defⁿ : The joint probability or cumulative distribution function of 2 r.v.s uniquely defines the probability of the joint events $\{X \leq x, Y \leq y\}$. It is defined by

$$F_{XY}(x, y) = P(X \leq x, Y \leq y) = \sum_{i=0}^{x_1} \sum_{j=0}^{y_1} P(x_i, y_j)$$

Properties:

$$\textcircled{1} \quad 0 \leq F_{XY}(x, y) \leq 1$$

$$\textcircled{2} \quad F_{XY}(-\infty, \infty) = 1.$$

$$\textcircled{3} \quad F_{XY}(-\infty, y) = F_{XY}(x, \infty) = 0.$$

* $\textcircled{4}$ F_{XY} is non-decreasing .

$$\textcircled{5} \quad F_X(x) = F_{XY}(x, \infty) = P(X \leq x) = P(X \leq x, Y \leq \infty).$$

$$\textcircled{6} \quad F_Y(y) = F_{XY}(\infty, y) = P(Y \leq y) = P(X \leq \infty, Y \leq y).$$

Marginal Probability Distribution Function .

$F_{XY}(x, \infty)$, $F_{XY}(\infty, y)$ are called marginal P.D.F of X & Y respectively . It is given by

$$F_{XY}(x, \infty) = \sum_j P(x, y_j)$$

$$F_{XY}(\infty, y) = \sum_i P(x_i, y)$$

eg: Suppose 2 dice are rolled simultaneously, then probability

of getting 3 on first die is ?

$$P(X=3) = P_{XY}(3, \infty) = P(X=3, Y=1) + P(X=3, Y=2)$$

$$+ P(X=3, Y=3) + P(X=3, Y=4)$$

$$+ P(X=3, Y=5) + P(X=3, Y=6)$$

$$= \frac{1}{36} + \frac{1}{36} + \frac{1}{36} + \frac{1}{36} + \frac{1}{36} + \frac{1}{36}$$

$$= \frac{1}{6} \approx$$

Continuous Multiple Random Variables:

Joint Probability density function: It is denoted

by $f_{xy}(x,y)$, also $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{xy}(x,y) dx dy = 1$

$$f_{XY}(x, y) dx dy = P(x \leq X \leq (x+dx), y \leq Y \leq (y+dy))$$

Note : $f_{xy}(x,y) = \frac{\partial^2 F_{xy}(x,y)}{\partial x \partial y}$

Relation b/w IPF & JPDF.

Perspectives:

$$\text{Properties: } \textcircled{1} \quad P(a \leq x \leq b, c \leq y \leq d) = \int_a^b \int_c^d f_{xy}(x, y) dx dy$$

$$\textcircled{2} \quad F_{XY}(x_1, y_1) = \int_{-\infty}^{x_1} \int_{-\infty}^{y_1} f_{XY}(x, y) dx dy$$

Marginal Probability density function (MPDF)

Marginal Probability Density Function (MPDF)

The MPDF of one of the r.v.s is the integral of joint p.d.f over the other r.v.

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x,y) dy \quad \text{is marginal p.d.f of } X.$$

$$f_y(y) = \int_{-\infty}^{\infty} f_{xy}(x,y) dx \quad \text{is} \quad n \quad n \quad n \quad y.$$

Statistical properties of Jointly Distributed Random Variables.

Defⁿ: 2 jointly distributed r.v.s X, Y are said to be statistically independent of each other iff

$$f_{XY}(x,y) = f_X(x)f_Y(y) \quad \forall (x,y) \text{ in the given range}$$

i.e joint p.d.f = product of 2 marginal p.d.f's.

Also if $F_{X_1X_2}(x_1, x_2) = F_X(x_1)F_X(x_2)$ & x_1, x_2 then also the 2 r.v.s X_1, X_2 are said to be statistically independent.

Conditional Probability distribution Density function.

Defⁿ: $f_Y(Y/x) = \frac{f_{XY}(x,y)}{f_X(x)}$ is conditional p.d.f of Y .

$f_X(X/Y) = \frac{f_{XY}(x,y)}{f_Y(y)}$ is conditional p.d.f of X .

Note: If X, Y are statistically independent then

$$\begin{aligned} f_X(X/Y) &= f_X(x) \\ &\& f_X(x/y) = \frac{f_{XY}(x,y)}{f_Y(y)} \quad \text{(by defⁿ)} \\ &\& f_Y(Y/x) = f_Y(y) \quad \left\{ \begin{array}{l} \therefore f_X(x/y) = \frac{f_{XY}(x,y)}{f_Y(y)} \\ \therefore \text{statistically independent} \\ \therefore f_X(x) \end{array} \right. \\ &\& \end{aligned}$$

Conditional Probability Distribution Function.

Let X, Y be 2 r.v.s discrete or continuous.

Defⁿ: Let X, Y be 2 r.v.s discrete or continuous. Then $F_Y(Y/x) = \frac{F_{XY}(x,y)}{F_X(x)}$, $F_X(x) > 0$ is conditional distribution of r.v. Y given $X=x$.

Also $F_X(X/Y) = \frac{F_{XY}(x,y)}{F_Y(y)}$, $F_Y(y) > 0$ is conditional distribution of r.v. X given $Y=y$.

Problems related to Binomial Distribution.

- ① Pg 125 Pb 3. A die is thrown 6 times. If getting an even no. is a success, find the probabilities of (i) atleast one success (ii) ≤ 3 successes (iii) 4 successes.

$$\text{W.R.T } P(X=r) = {}^n C_r p^r q^{n-r}.$$

Sol: Given $n=6$.

p = getting an even no. is success.
we have 2, 4, 6 as even nos. on die.
 $\therefore p = \frac{3}{6} = \frac{1}{2} \Rightarrow q = \frac{1}{2}$ (\because Total faces of dice = 6
Even " " = 3)

$$(i) P(X \geq 1) = 1 - P(X=0) \\ = 1 - {}^6 C_0 p^0 q^{6-0} = 1 - q^6 = 1 - \left(\frac{1}{2}\right)^6 = \frac{63}{64}$$

$$(ii) P(X \leq 3) = P(X=0) + P(X=1) + P(X=2) + P(X=3) \\ = {}^6 C_0 \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^6 + {}^6 C_1 \left(\frac{1}{2}\right)^1 \left(\frac{1}{2}\right)^5 + {}^6 C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^4 + {}^6 C_3 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^3 \\ = \left(\frac{1}{2}\right)^6 \left[1 + 6 + \frac{6 \times 5}{2} + \frac{6 \times 5 \times 4}{2 \times 3} \right] \\ = \frac{1}{64} [1 + 6 + 15 + 20] = \frac{42}{64} = \frac{21}{32}.$$

$$(iii) P(X=4) = {}^6 C_4 p^4 q^{6-4} = \frac{3 \times 5}{2} \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^2 = \frac{15}{64}$$

- ② Pg 126 Pb 4: 10 coins are thrown simultaneously. Find the probability of getting atleast (i) 7 heads (ii) 6 heads.

Sol: Here $n=10$, p = getting head = $\frac{1}{2} \Rightarrow q = \frac{1}{2}$

$$(i) P(X \geq 7) = P(X=7) + P(X=8) + P(X=9) + P(X=10) \\ = {}^{10} C_7 \left(\frac{1}{2}\right)^7 \left(\frac{1}{2}\right)^3 + {}^{10} C_8 \left(\frac{1}{2}\right)^8 \left(\frac{1}{2}\right)^2 + {}^{10} C_9 \left(\frac{1}{2}\right)^9 \left(\frac{1}{2}\right)^1 + {}^{10} C_{10} \left(\frac{1}{2}\right)^{10} \\ = \binom{10}{7} \left[\frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2} + \frac{10 \times 9}{2} + 10 + 1 \right] \\ = \frac{(120+45+1)}{1024} = \frac{176}{1024} = \frac{11}{64}.$$

$$(ii) P(X \geq 6) = P(X=6) + P(X \geq 7) \\ = {}^{10} C_6 \left(\frac{1}{2}\right)^6 \left(\frac{1}{2}\right)^4 + \frac{11}{64} = \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2 \times 1} \cdot \frac{1}{1024} + \frac{176}{1024} \\ = \frac{210+176}{1024} = \frac{386}{1024} = \frac{193}{512}.$$

- ③ If 3 of 20 tyres are defective & 4 of them are randomly chosen for inspection, what is the probability that only one of the defective tyre will be included? (Pg 127 Pb 7)

Sol: Given = 20 tyres; 3 are defective \Rightarrow 17 are good $\Rightarrow p = \frac{3}{20}$ is defective tyres probability.
To chose 4 at random $\Rightarrow n=4$.

$$P(X=1) = ? \quad {}^nC_1 p^1 q^{n-1} \\ = {}^4C_1 \left(\frac{3}{20}\right) \left(\frac{17}{20}\right)^3 = \frac{4 \times 3 \times 17^3}{(20)^4} = 0.3685.$$

\rightarrow Pg 127 Pb 7.

- ④ Determine B.D for which mean is 4 & variance 3 (Pg 128 Pb 9)

$$B.D = b(x; n, p) \\ B.D = b(x; n, p) \quad \text{Give } \mu = 4, \sigma^2 = 3 \\ \text{To find } n; p \\ \text{w.r.t } \mu = np \Rightarrow np = 4 \quad \text{Also } \sigma^2 = npq = 3 \Rightarrow npq = 3 \\ \therefore npq = 3 \Rightarrow 4q = 3 \Rightarrow q = \frac{3}{4} \Rightarrow p = 1-q = \frac{1}{4} \\ \therefore n\left(\frac{1}{4}\right) = 4 \Rightarrow n = 16.$$

$$\text{Thus } b(x; n, p) = b(x; 16, \frac{1}{4}).$$

- ⑤ In 256 sets of 12 tosses of a coin, in how many cases one can expect 8 heads & 4 tails. (Pg 128 Pb 11)

$$\text{Sol: } N = 256, \quad n = 12, \quad p = \text{getting head or success} = \frac{1}{2} \Rightarrow q = \frac{1}{2} \\ P(X=8) = ? = {}^{12}C_8 p^8 q^4 = {}^{12}C_8 \left(\frac{1}{2}\right)^8 \left(\frac{1}{2}\right)^4 = \frac{12 \times 11 \times 10 \times 9}{4 \times 3 \times 2} \left(\frac{1}{2}\right)^{12} \\ = \frac{495}{(2)^{12}}$$

$$\therefore \text{Expected no. of such cases in 256 sets} \\ = N \cdot P(X=8) = \frac{(256)(495)}{(2)^{12}} = \frac{8}{2^4} \frac{(495)}{2^8} = \frac{495}{2^4}$$

$$= \frac{495}{16}$$

$$= 30.9375$$

$$\approx 31$$

6. Six dice are thrown 729 times. How many times do you expect at least three dice to show a 5 or 6? (Pg 129 Pb 13)

$$\text{Sol: } N = 729, m = 6, p = \text{to show 5 or 6} = \frac{2}{6} = \frac{1}{3} \Rightarrow q = \frac{2}{3}$$

$$\begin{aligned} P(X \geq 3) &=? = {}^nC_3 p^3 q^{n-3} + {}^nC_4 p^4 q^{n-4} + {}^nC_5 p^5 q^{n-5} + {}^nC_6 p^6 q^{n-6} \\ &= P(X=3) + P(X=4) + P(X=5) + P(X=6) \\ &= {}^6C_3 \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^3 + {}^6C_4 \left(\frac{1}{3}\right)^4 \left(\frac{2}{3}\right)^4 + {}^6C_5 \left(\frac{1}{3}\right)^5 \left(\frac{2}{3}\right)^5 + {}^6C_6 \left(\frac{1}{3}\right)^6 \\ &= \frac{4 \times 5 \times 4}{3 \times 2} \times \frac{8}{3^6} + \frac{6 \times 5}{2} \times \frac{16}{3^6} + 6 \times \frac{2}{3^6} + \frac{1}{3^6} \\ &= \frac{160 + 60 + 12 + 1}{3^6} = \frac{233}{3^6} = \frac{233}{729} \end{aligned}$$

∴ Expected no. of such cases in 729 times = $N \cdot P(X \geq 3)$

$$= 229 \left(\frac{233}{729}\right) = \underline{\underline{233}}$$

* 7 (i) Out of 800 families with 5 children each, how many would you expect to have (a) 3 boys (b) 5 girls (c) either 2 or 3 boys (d) atleast one boy? Assume equal probabilities for boys & girls.

(ii) Out of 800 families with 4 children each, how many families would be expected to have (a) 2 boys & 2 girls (b) atleast one boy (c) no girl (d) atleast one girl? Assume equal probabilities for boys & girls.
(Pg 131 Pb 19)

$$\text{Sol: (i) } N = 800, m = 5, p = \text{probability of child to be boy} = \frac{1}{2} \Rightarrow q = \frac{1}{2}$$

$$(a) P(X=3) = ? = {}^5C_3 p^3 q^2 = {}^5C_3 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^2 = \frac{5 \times 4}{2} \times \frac{1}{4 \times 8} = \frac{5}{16}$$

$$(b) P(X=0) = {}^5C_0 p^0 q^5 = \left(\frac{1}{2}\right)^5 = \frac{1}{32}$$

$$(c) P(X=2) + P(X=3) = {}^5C_2 p^2 q^3 + {}^5C_3 p^3 q^2 = \frac{5 \times 4}{2} \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^3 + \frac{5 \times 4}{2} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^2 = \frac{20}{8} = \frac{5}{8}$$

$$(d) P(X \geq 1) = 1 - P(X=0) = 1 - \frac{1}{32} = \frac{31}{32}$$

Thus (a) For 800 families expectation of 3 boys = $800 \cdot P(X=3)$

$$\Rightarrow 800 \times \frac{5}{16} = \underline{\underline{250}} \text{ families}$$

(b) Expected no. of families to have 5 girls = $N \cdot P(x=0)$
 $= \frac{25}{800} \times \frac{1}{32} = \underline{25}$ families.

(c) Expected no. of families to have either 2 or 3 boys
 $= N \cdot [P(x=2) + P(x=3)] = \frac{100}{800} \times \frac{5}{8} = \underline{50}$ families.

(d) Expected no. of families with atleast one boy
 $= N \cdot P(x \geq 1) = \frac{25}{800} \times \frac{31}{32} = \underline{775}$ families.

(i) $N = 800, n=4, p = \text{probability of child to be boy} = \frac{1}{2} \Rightarrow q = \frac{1}{2}$.

(a) $P(x=2) = {}^nC_2 p^2 q^2 = {}^4C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^2 = \frac{4 \times 3}{2} \times \frac{1}{4 \times 4} = \frac{3}{8}$

\therefore Expected no. of families to have 2 boys & 2 girls
 $= N \cdot P(x=2) = \frac{100}{800} \times \frac{3}{8} = \underline{37.5}$ families.

(b) $P(x \geq 1) = 1 - P(x=0) = 1 - {}^4C_0 p^0 q^4 = 1 - \left(\frac{1}{2}\right)^4 = \frac{15}{16}$

\therefore Expected no. of families to have atleast 1 boy
 $= N \cdot P(x \geq 1) = \frac{50}{800} \times \frac{15}{16} = \underline{750}$ families.

(c) $P(x=4) = {}^4C_4 \left(\frac{1}{2}\right)^4 = \frac{1}{16} \quad \left\{ \because \text{No girl} \Rightarrow \text{all 4 are boys} \right\}$

\therefore Expected no. of families to have ~~atleast~~ ^{NO} girls
 $= N \cdot P(x=4) = \frac{50}{800} \times \frac{1}{16} = \underline{50}$ families.

(d) $P(x=4) + P(x=3) \quad \left(\begin{array}{l} \text{Atleast 1 girl} \Rightarrow \text{No girl \& 1 girl} \\ \Rightarrow 4 \text{ Boys \& 3 boys} \end{array} \right)$
 $= {}^4C_4 \left(\frac{1}{2}\right)^4 + {}^4C_3 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)$

$$= \frac{1}{16} + 4 \times \frac{1}{16} = \frac{5}{16}$$

\therefore Expected no. of families to have atleast 1 girl

$$= N \cdot [P(x=3) + P(x=4)]$$

$$= \frac{50}{800} \times \frac{5}{16} =$$

$$= \underline{250} \text{ families}$$

⑧ (Pg 135 Pb25) In a B.D consisting of 5 independent trials, probabilities of 1 & 2 success are 0.4096 & 0.2048 respectively. Find the parameter 'p' of the distribution.

Sol: $n=5$, $P(X=1)=0.4096$; $P(X=2)=0.2048$. $p=?$

$$P(X) = {}^n C_r P^r q^{n-r}$$

$$\therefore P(X=1) = {}^5 C_1 p q^{5-1} = 0.4096 \Rightarrow 5 p q^4 = 0.4096 \quad (1)$$

$$P(X=2) = {}^5 C_2 p^2 q^{5-2} = 0.2048 \Rightarrow \frac{5 \times 4}{2} p^2 q^3 = 10 p q^3 = 0.2048 \quad (2)$$

$$\frac{(1)}{(2)} \Rightarrow \frac{P(X=1)}{P(X=2)} = \frac{5 p q^4}{10 p q^3} = \frac{0.4096}{0.2048}$$

$$\Rightarrow \frac{q}{2p} = 2 \Rightarrow q = 4p \quad \left\{ \text{as } p = 1 - q \right\} \Rightarrow 1 - p = 4p$$

OR

By recurrence relation $P(r+1) = \frac{(n-r)p}{(r+1)q} P(r)$.

$$\Rightarrow 1 = 5p \quad \Rightarrow p = \frac{1}{5}$$

$$\Rightarrow p = 0.2$$

$$\begin{aligned} \text{Thus } P(X=2) &= \frac{(n-r)p}{(r+1)q} P(X=1) \\ &= \frac{(5-1)p}{2q} P(X=1) \\ \Rightarrow 0.2048 &= \frac{4p}{2q} (0.4096) \\ \Rightarrow q &= 4p \Rightarrow 1 - p = 4p \Rightarrow 1 = 5p \Rightarrow p = \frac{1}{5} = 0.2 \end{aligned}$$

⑨ (Pg 138 Pb32) A coin is biased in a way that a head is twice as likely to occur as likely to occur as a tail. If the coin is tossed 3 times, find the probability of getting 2 tail & 1 head.

Sol: Let p = probability of getting tail as success. $\Rightarrow n=3$.

Given $P(H) = 2 P(T)$ w.r.t $P(H) + P(T) = 1 \Rightarrow 3 P(T) = 1 \Rightarrow P(T) = \frac{1}{3} \Rightarrow p = \frac{1}{3}$.

Given $P(H) = 2 P(T)$ w.r.t $P(H) + P(T) = 1 \Rightarrow 3 P(T) = 1 \Rightarrow P(T) = \frac{1}{3} \Rightarrow p = \frac{1}{3}$.

To find $P(2T) = ?$. $P(X=2) = {}^n C_r p^r q^{n-r}$

$$\begin{aligned} &= {}^3 C_2 \left(\frac{1}{3}\right)^2 \left(\frac{2}{3}\right)^{3-2} \\ &= 3 \left(\frac{1}{3}\right)^2 \times \frac{2}{3} = \underline{\underline{\frac{2}{9}}} \end{aligned}$$

(10) Pg 139 Pb 33. Fit a B.D to the following frequency distribution.

x	0	1	2	3	4	5	6
f	13	25	52	58	32	16	4

Fit a B.D \Rightarrow Obtain $E(x)$ where $E(x) = N \cdot p(x)$; $N = \sum_i f_i$
 $\& p(x) = f(x)$.

From data $n = 6$ ($\because x = 0, 1, \dots, 6$) ; $N = \sum_i f_i = 13 + 25 + 52 + 58 + 32 + 6 + 4$

$$\text{Mean} = \frac{\sum x_i f_i}{\sum f_i} = np \Rightarrow \frac{(0 \times 13) + (1 \times 25) + (2 \times 52) + (3 \times 58) + (4 \times 32) + (5 \times 16)}{200} = 2.00$$

$$\Rightarrow np = \frac{535}{200} = 2.675$$

$$\therefore n=6 \Rightarrow np=2.675 \Rightarrow p = 0.446 \Rightarrow q = 0.554 .$$

$$P(x=0) = {}^n C_0 p^0 q^{n-0} = {}^6 C_0 (0.554)^6 = (0.554)^6 = 0.02891$$

$$P(x=1) = {}^6 C_1 (0.554)^5 (0.446)^1 = 6 (0.446)(0.554)^5 = 0.1396$$

$$P(x=2) = {}^6 C_2 (0.446)^2 (0.554)^4 = \frac{15}{2} (0.446)^2 (0.554)^4 = 0.2809$$

$$P(x=3) = {}^6 C_3 (0.446)^3 (0.554)^3 = \frac{20}{3} (0.446)^3 (0.554)^3 = 0.3016$$

$$P(x=4) = {}^6 C_4 (0.446)^4 (0.554)^2 = \frac{15}{2} (0.446)^4 (0.554)^2 = 0.1821$$

$$P(x=5) = {}^6 C_5 (0.446)^5 (0.554)^1 = 6 (0.446)^5 (0.554) = 0.05864$$

$$P(x=6) = {}^6 C_6 (0.446)^6 (0.554)^0 = (0.446)^6 = 0.007866$$

\therefore B.D is

x	0	1	2	3	4	5	6
f	13	25	52	58	32	16	4
$E(x)$	6	28	56	60	36	12	2

where $E(x)$ values are rounded off to nearest integer.

(11) (Pg 147 Pb 45) Find the probability of getting an even number 3 or 4 or 5 times in throwing 10 dice, using B.D.

Sol: Let P = prob. of getting even no. We have 2, 4, 6 as possible even nos. on a die. $\therefore P = \frac{3}{6} = \frac{1}{2} \Rightarrow q = \frac{1}{2}$.

Given $n = 10$. To find $P(x=3) = ?$, $P(x=4) = ?$, $P(x=5) = ?$

$$\therefore P(x=3) = {}^{10} C_3 \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^7 = \frac{10 \times 9 \times 8}{3 \times 2 \times 1} \left(\frac{1}{2}\right)^{10} = \frac{120}{1024} = 0.112 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\}$$

$$P(x=4) = {}^{10} C_4 \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^6 = \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2 \times 1} \left(\frac{1}{2}\right)^{10} = \frac{210}{1024} = 0.2 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} .558$$

$$P(x=5) = {}^{10} C_5 \left(\frac{1}{2}\right)^5 \left(\frac{1}{2}\right)^5 = \frac{10 \times 9 \times 8 \times 7 \times 6}{5 \times 4 \times 3 \times 2 \times 1} \left(\frac{1}{2}\right)^{10} = \frac{252}{1024} = 0.246 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\}$$

- ⑥ Find the probability of getting an even number 3 or 4 or 5 times with in throwing 10 dice, using B.D.

Sol: Given $n = 10$
 We know no. of even nos. on dice = 3 ($\because 2, 4, 6$ are the nos.)
 $\therefore p = \text{prob. of getting an even number} = \frac{3}{6} = \frac{1}{2}$
 $\therefore q = 1 - p = \frac{1}{2}$
 $P(X=3) + P(X=4) + P(X=5) = {}^{10}C_3 p^3 q^7 + {}^{10}C_4 p^4 q^6 + {}^{10}C_5 p^5 q^5$
 $= \frac{10 \times 9 \times 8}{3 \times 2} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^7 + \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2} \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^6 + \frac{10 \times 9 \times 8 \times 7 \times 6}{5 \times 4 \times 3 \times 2} \left(\frac{1}{2}\right)^5 \left(\frac{1}{2}\right)^5$
 $= \left(\frac{1}{2}\right)^{10} [120 + 210 + 252] = 0.112 + 0.2 + 0.246$
 $= \frac{582}{1024} =$

- ⑦ A coin is biased in a way that a head is twice as likely to occur as a tail. If the coin is tossed 3 times, find the probability of getting 2 tail and 1 head.

Sol: $P(H) = 2 P(T)$; Given $n=3$, To find $P(X=2 \text{ tail } \& 1 \text{ Head})=?$
 W.K.t. $P(H) + P(T) = 1 \Rightarrow 2P(T) + P(T) = 1 \Rightarrow P(T) = \frac{1}{3}$
 $\therefore P(H) = \frac{2}{3}$.
 $P(X=2) = ? = {}^3C_2 p^2 q^1 = \frac{3 \times 2}{2} \left(\frac{2}{3}\right)^2 \left(\frac{1}{3}\right) = \underline{\underline{\frac{2}{9}}} \text{ Ans.}$

- ⑧ Out of 800 families with 5 children each, how many could you expect to have
 a) 3 boys b) 5 girls c) either 2 or 3 boys d) atleast 1 boy?
 Assume equal probabilities for boys & girls.

Sol: Given $N = 800$, $n = 5$. Let p - prob. of getting boy is success.
 $\therefore p = \frac{1}{2}, q = \frac{1}{2}$ { \because equal probabilities for boys & girls}.

a) $P(X=3) = {}^5C_3 p^3 q^2 = \frac{5 \times 4}{2} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^2 = \frac{5}{8} \text{ per family.}$
 $\therefore \text{Total no. of families with 3 boys} = N \cdot P(x) = \frac{800 \times 5}{8} = \underline{\underline{500}} \text{ families}$

b) $P(X=0) (\because \text{all are girls i.e. 5 girls}) = {}^5C_0 p^0 q^5 = \left(\frac{1}{2}\right)^5 = \frac{1}{32}$
 $\therefore \text{Total no. of families with 5 girls} = \frac{2500}{32} = 25 \text{ families.}$

$$\textcircled{c} \quad P(X=2) + P(X=3) = {}^5C_2 p^2 q^3 + {}^5C_3 p^3 q^2 = \frac{2}{5} \cdot \frac{5 \times 4}{2} \left(\frac{1}{2}\right)^5 + \frac{3}{5} \cdot \frac{4 \times 3 \times 2}{3} \left(\frac{1}{2}\right)^5 = \frac{20}{32} \cdot \frac{5}{8} = \frac{5}{8}$$

∴ total no. of families with either 2 or 3 boys = $\frac{800 \times 5}{8} = 500$ families.

\textcircled{d}

Problems on Poisson Distribution.

- (1) (Pg 159 Pb 5) If a bank received on the average 6 bad cheques per day, find the probability that it will receive 4 bad cheques on any day. Given $\lambda = 6$. $P(X=4) = 0.1339$.

Sol: $P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$; Given $\lambda = 6$. To find $P(X=4)$?

$$\therefore P(X=4) = \frac{e^{-6} 6^4}{4!} = \frac{54}{e^6} = 0.1339$$

- (2) A car-hire firm has cars which it hires out day by day. The no. of demands for a car on each day is distributed as a Poisson distribution with mean 1.5. Calculate the proportion of days (i) on which there is no demand (ii) on which demand is refused. (Pg 158 Pb 2)

Sol: Given $\lambda = 1.5$: Total cars = 2; $N = 365$ days in a year
 No demand \Rightarrow Cars were not hired $\Rightarrow P(X=0) = ?$
 No. of days when there is no demand in a year $= N \cdot P(X=0)$

$$(i) P(X=0) = e^{-1.5} \frac{(1.5)^0}{0!} = \frac{e^{-1.5}}{1} = 0.2231$$

\therefore No. of days in a year when there is no demand of car $= N \cdot P(X=0)$
 $= 365 \times 0.2231$
 $= 81$ days

- (ii) First we find probability on which demand is refused. \Rightarrow refusal takes place only when demand is for more than 2 cars. i.e. to find $P(X > 2)$

To calculate proportion of days when there is no refusal if demand $= E(X > 2) = ?$

$$P(X > 2) = 1 - [P(X=0) + P(X=1) + P(X=2)] = 1 - e^{-1.5} [1 + 1.5 + \frac{1.5^2}{2}]$$

$$= 1 - e^{1.5} [0 + (1.5) + \frac{(1.5)^2}{2}] = 0.1913$$

$$E(X > 2) = N P(X > 2) = 365 (0.1913) = 69.82 \approx \underline{\underline{70 \text{ days}}}$$

- (3) (Pg 160 P&P) It has been found that 2% of the tools produced by a certain machine are defective. What is the probability that in a shipment of 400 such tools (a) 3% or more
(b) 2% or less will prove defective.

Sol: No. of tools = 400.

$$\% \text{ of defective tools} = 2\% = p$$

$$\lambda = \text{No. of defective tools for 400} = np = (2\%) \times (400) \\ = 400 \times \frac{2}{100} \\ \Rightarrow \boxed{\lambda = 8}$$

$$(i) P(X \geq 3\%) = ?$$

$$3\% \Rightarrow \frac{3}{100} \times 400 = 12 \text{ tools}$$

$$\therefore P(X \geq 3\%) = P(X \geq 12) = 1 - P(X \leq 11)$$

$$= 1 - e^{-\lambda} \left[\frac{\lambda^0}{0!} + \frac{\lambda^1}{1!} + \frac{\lambda^2}{2!} + \dots + \frac{\lambda^{11}}{11!} \right]$$

$$= 1 - e^{-8} \left[1 + 8 + \frac{8^2}{2!} + \frac{8^3}{3!} + \dots + \frac{8^{11}}{11!} \right]$$

$$= 1 - e^{-8} [2647.29]$$

$$= 0.1119.$$

$$(ii) P(X \leq 2\%)$$

$$2\% \Rightarrow \frac{2}{100} \times 400 = 8 \text{ tools}$$

$$\therefore P(X \leq 2\%) = P(X \leq 8) = P(X=0) + P(X=1) + P(X=2) + P(X=3) + P(X=4) \\ + P(X=5) + P(X=6) + P(X=7) + P(X=8)$$

$$= e^{-\lambda} \left[\frac{\lambda^0}{0!} + \frac{\lambda^1}{1!} + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \frac{\lambda^4}{4!} + \frac{\lambda^5}{5!} + \frac{\lambda^6}{6!} + \frac{\lambda^7}{7!} + \frac{\lambda^8}{8!} \right]$$

$$= e^{-8} \left[1 + 8 + \frac{8^2}{2!} + \frac{8^3}{3!} + \frac{8^4}{4!} + \frac{8^5}{5!} + \frac{8^6}{6!} + \frac{8^7}{7!} + \frac{8^8}{8!} \right] \frac{16777216}{40320}$$

$$= e^{-8} [9 + 32 + 85.3333 + 170.6667 + 273.0667 + 3,640.89 \\ + 416.1016 + 416.1016]$$

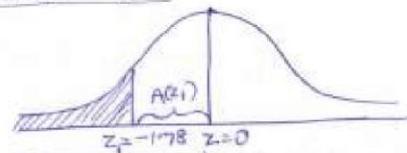
$$= e^{-8} \boxed{475642.3588} = \underline{\underline{0.5995}}$$

Problems related to Normal Distribution.

1. (Pg 205 pb 10) If X is a normal variate, find the area A
- to the left of $z = -1.78$
 - to the right of $z = -1.45$
 - corresponding to $-0.8 \leq z \leq 1.53$
 - to the left of $z = -2.52$ and to the right of $z = 1.83$.

Sol: (i) Required Area ' A ' is 'shaded region'.

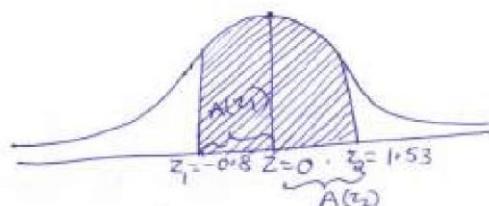
$$\begin{aligned} \text{Area} &= 0.5 - A(z_1) \\ &= 0.5 - A(-1.78) \\ &= 0.5 - A(1.78) \\ &= 0.5 - 0.4625 \quad (\text{from tables}) \\ &= 0.0375 \end{aligned}$$



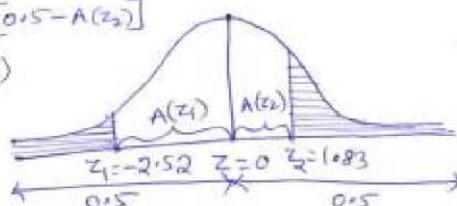
$$\begin{aligned} \text{(ii) Required Area} &= 0.5 + A(z_1) \\ &= 0.5 + A(-1.45) \\ &= 0.5 + A(1.45) \quad (\text{By symmetry}) \\ &= 0.5 + 0.4265 \\ &= 0.9265 \end{aligned}$$



$$\begin{aligned} \text{(iii) Area} &= A(z_1) + A(z_2) \\ &= A(0.8) + A(1.53) \\ &= A(0.8) + A(1.53) \\ &= 0.2881 + 0.437 \\ &= 0.7251 \end{aligned}$$



$$\begin{aligned} \text{(iv) Required Area} &= [0.5 - A(z_1)] + [0.5 - A(z_2)] \\ &= 0.5 - A(-2.52) + 0.5 - A(1.83) \\ &= 1 - A(2.52) - A(1.83) \\ &= 1 - 0.4941 - 0.4664 \\ &= 1 - 0.9605 \\ &= 0.0395 \end{aligned}$$



2. (Pg 208 Pb13) If the masses of 300 students are normally distributed with mean 68 kgs & standard deviation 3 kgs, how many students have masses

- (i) Greater than 72 kgs
- (ii) Less than or equal to 64 kgs
- (iii) Between 65 & 71 kgs inclusive.

Sol. No. of students = 300 = N.

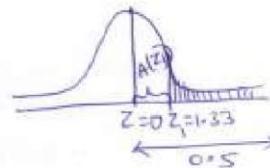
$$\mu = 68 ; \sigma = 3 .$$

(i) No. of students with mass greater than 72 kgs = $N \cdot P(X)$.

when $X=72$

$$P(X > 72) = P\left(Z > \frac{72-68}{3}\right) = P(Z > 1.33)$$

$$\begin{aligned} \therefore P(X > 72) &= P(Z > 1.33) \\ &= 0.5 - A(2) \\ &= 0.5 - A(1.33) \\ &= 0.5 - 0.4082 \end{aligned}$$

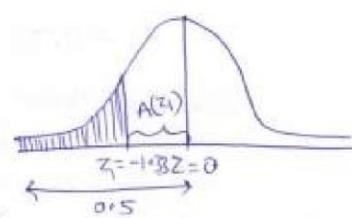


$$\therefore \text{No. of students} = 300 \times 0.0918 = 27 \text{ students}.$$

(ii) less than or equal to 64 kgs = $P(X \leq 64) = ?$

$$Z = \frac{X-\mu}{\sigma} = \frac{64-68}{3} = -\frac{4}{3} \approx -1.33$$

$$\begin{aligned} \therefore P(X \leq 64) &= P(Z \leq -1.33) \\ &= 0.5 - A(2) \\ &= 0.5 - A(-1.33) \\ &= 0.5 - A(1.33) \\ &= 0.5 - 0.4082 = 0.0918 . \end{aligned}$$

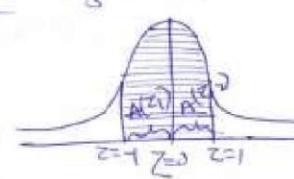


$$\therefore \text{No. of Students} = 300 \times 0.0918 = 27 \text{ students}.$$

$$(iii) P(65 \leq X \leq 71) \quad Z_1 = \frac{65-68}{3} = -1 ; Z_2 = \frac{71-68}{3} = 1$$

$$\begin{aligned} &= P(-1 \leq Z \leq 1) \\ &= A(2) + A(1) \\ &= 2A(1) \\ &= 2(0.3413) \\ &= 0.6828 \end{aligned}$$

$\therefore \text{No. of Students}$
$= 300 \times 0.6828$
$= 205 \text{ students}$



3. (Pg 213 Pb 20) If x is normally distributed with mean 2 & S.D. 0.1, then find $P(|x-2| > 0.01)$?

$$\text{Sol: } \mu = 2, \sigma = 0.1$$

$$P(|x-2| > 0.01) = 1 - P(|x-2| < 0.01)$$

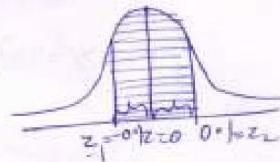
$$|x-2| < 0.01 \Rightarrow x = 2+0.01 \quad \& x = +2\pm 0.01 \\ = 2.01 \quad \& +1.99$$

$$\text{Let } x_1 = 1.99, x_2 = 2.01$$

$$z = \frac{x-\mu}{\sigma} \Rightarrow z_1 = \frac{1.99-2}{0.1} = \frac{-0.01}{0.1} = -0.1$$

$$\text{Also } z_2 = \frac{2.01-2}{0.1} = \frac{0.01}{0.1} = 0.1$$

$$\begin{aligned} \therefore P(|x-2| > 0.01) &= 1 - P(z_1 < z < z_2) \\ &= 1 - P(-0.1 < z < 0.1) \\ &= 1 - [A(z_1) + A(z_2)] \\ &= 1 - 2A(0.1) \\ &= 1 - 2(0.0398) \\ &= 1 - 0.0796 = 0.9204 \end{aligned}$$



4. (Pg 222 Pb 20) Suppose the weights of 800 male students are normally distributed with mean 28.8 kg & S.D. 2.06 kg. Find the no. of students whose weights are b/w (i) 28.4 kg & 30.4 kg (ii) more than 31.3 kg.

$$N = 800, \text{ mean } \mu = 28.8 \text{ ; } \sigma = 2.06$$

$$(i) P(28.4 \leq x \leq 30.4) = ?$$

$$\text{No. of Students} = N \times P(28.4 \leq x \leq 30.4)$$

$$; z = \frac{x-\mu}{\sigma}$$

$$z_1 = \frac{28.4 - 28.8}{2.06} = \frac{-0.4}{2.06} = -0.1942 = 0$$

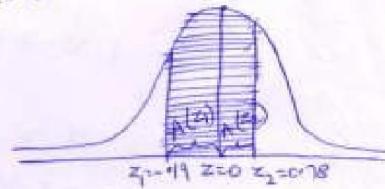
$$z_2 = \frac{30.4 - 28.8}{2.06} = \frac{1.6}{2.06} = \frac{80}{103} = 0.7761$$

$$z_2 = \frac{30.4 - 28.8}{2.06} = \frac{1.6}{2.06} = \frac{80}{103} = 0.7761$$

$$= \frac{1.6}{2.06} = \frac{80}{103} = 0.7761$$

$$\begin{aligned}
 P(28.4 \leq X \leq 30.4) &= P(z_1 \leq Z \leq z_2) \\
 &= P(-0.19 \leq Z \leq 0.78) \\
 &= A(z_1) + A(z_2) \\
 &= A(-0.19) + A(0.78) \\
 &= 0.0753 + 0.2823 \\
 &= 0.3576
 \end{aligned}$$

\therefore No. of students = $N P(x \geq 28.4)$
 $= 800 \times 0.3576$
 $= 286$

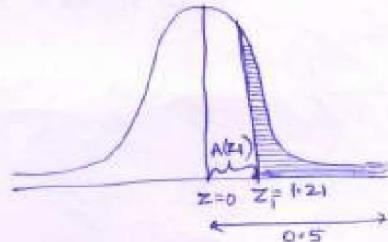


(ii) more than 31.3 kg.

$$P(X \geq 31.3)$$

$$z = \frac{31.3 - 28.8}{2.06} = \frac{2.5}{2.06} = 1.21$$

$$\begin{aligned}
 P(X \geq 31.3) &= P(Z \geq 1.21) \\
 &= 0.5 - A(z_1) \\
 &= 0.5 - A(1.21) \\
 &= 0.5 - 0.3869
 \end{aligned}$$



$$= 0.1131$$

(Q200 Pg3) \therefore No. of students = $N P(X \geq 31.3) = 800 \times 0.1131 = 90$
 Q5 In a Normal distribution, 7% of the items are under 35 and 89% are under 63. Determine mean & variance of the distribution.

$$P(X < 35) = 7\% \quad P(X < 63) = 89\% \quad \mu = ? \quad \sigma^2 = ?$$

$$x_1 = 35, \quad x_2 = 63 \Rightarrow z_1 = \frac{35-\mu}{\sigma}, \quad z_2 = \frac{63-\mu}{\sigma}$$

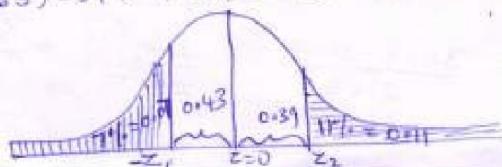
"Total 100% we consider values of $x < 50$ to be left of $Z = 0$ i.e. $x = \mu$ & values of $x > 50$ to be right of $Z = 0$. $\Rightarrow z_1 < 0$ & $z_2 > 0 \Rightarrow \frac{35-\mu}{\sigma} = z_1$ (say)

$$\frac{63-\mu}{\sigma} = z_2$$

$$P(X < 35) = 7\% \Rightarrow A(Z < z_1) = 7\% \Rightarrow A(z_1) = 0.5 - 0.07 = 0.43$$

$$P(X < 63) = 89\% \Rightarrow A(Z < z_2) = 89\% \Rightarrow A(z_2) = 0.5 - 0.01 = 0.49$$

$$\Rightarrow A(z_2) = 0.5 - 0.11 = 0.39$$



Next page
Page

Given $P(X < 35) = 7\% = 0.07$; $P(X < 63) = 89\% = 0.89$

when $x = 35$, $z = \frac{x-\mu}{\sigma} = \frac{35-\mu}{\sigma} = -z_1$ say $\rightarrow \textcircled{1}$

When $x = 63$, $z = \frac{x-\mu}{\sigma} = \frac{63-\mu}{\sigma} = z_2$ say $\rightarrow \textcircled{2}$

From the fig., $P(0 < z < z_2) = 0.39 \Rightarrow A(z_2) = 0.39$
 $\& P(0 < z < z_1) = 0.43 \Rightarrow A(z_1) = 0.43$

Using tables we find $z_1 = 1.48$, $z_2 = 1.23$
 $\Rightarrow z_1 = -1.48$, $z_2 = 1.23$

$$\therefore \textcircled{1} \& \textcircled{2} \text{ are } \frac{35-\mu}{\sigma} = -1.48 \Rightarrow 35 - \mu = -1.48\sigma$$

$$\Rightarrow \mu + 1.48\sigma = 35 \quad \textcircled{3}$$

$$\frac{63-\mu}{\sigma} = 1.23 \Rightarrow 63 = \mu + 1.23\sigma \Rightarrow \mu + 1.23\sigma = 63 \quad \textcircled{4}$$

$$\text{Solving } \textcircled{3} \& \textcircled{4} \quad \sigma = \frac{28}{2.71} = 10.332$$

$$\Rightarrow \sigma^2 = 106.75$$

$$\text{From } \textcircled{3} \quad 35 = \mu - 1.48(10.332) \Rightarrow \mu = 35 + 15.3 = 50.3$$

$$\therefore \mu = 50.3$$

- * 6 In a ND 3.I. of the items are under 45 & 8.I. are over 64.
 Find the mean & variance of the distribution. (Pg 202 Pg 4)

$$P(X < 45) = 31\% = 0.31; P(X > 64) = 8\% = 0.08$$

$$z = \frac{x-\mu}{\sigma} \Rightarrow z_1 = \frac{45-\mu}{\sigma}; z_2 = \frac{64-\mu}{\sigma}$$

The fig. is:

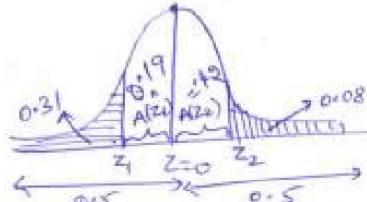
From table

$A(z = 0.19)$ when $z = 0.5$

$$\Rightarrow z_1 = -0.5$$

Area = 0.42 when $z = 1.4$

$$\Rightarrow z_2 = 1.4$$



$$\Rightarrow \frac{45-\mu}{\sigma} = -0.5 \Rightarrow 45-\mu = -0.5\sigma \Rightarrow \mu + 0.5\sigma = 45 \quad \text{---(1)}$$

$$\cdot \frac{64-\mu}{\sigma} = 1.4 \Rightarrow 64 = \mu + 1.4\sigma \quad \text{---(2)}$$

From (1) & (2) we get $\sigma = \frac{10}{1.4} = 10 \Rightarrow \sigma = 10$

From (1) : $\mu = 45 + 0.5\sigma = 45 + 0.5(10) = 50 \Rightarrow \mu = 50$

(Ques 203 P67) The marks obtained in mathematics by 1000 students is normally distributed with mean 78% & S.D 11%. Determine

- (i) How many students got marks above 90%.
- (ii) What was the highest mark obtained by the lowest 10% of the students.

(iii) Within what limits did the middle of 90% of the students lie.

Sol: No. of students = 1000 ; $\mu = 78\% = 0.78$; $\sigma = 11\% = 0.11$

(i) $P(X > 90\%) = ?$ $z = \frac{x-\mu}{\sigma} = \frac{0.9 - 0.78}{0.11} = 1.09 = z_1$

$$P(X > 90\%) = P(z > z_1) = P(z > 1.09) = 0.5 - A(z_1 = 1.09)$$

$$= 0.5 - 0.3621$$

$$= 0.1379$$

No. of students who got marks above 90% = $1000 P(X > 90\%)$
 $= 1000 \times 0.1379$
 $= 137.9$ approx
 $= 138$ students.

(ii) $P(X < x_1) = 10\% = 0.1$

$$\Rightarrow P(z < z_1) = 0.1$$

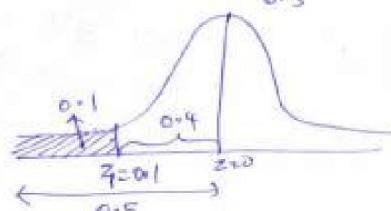
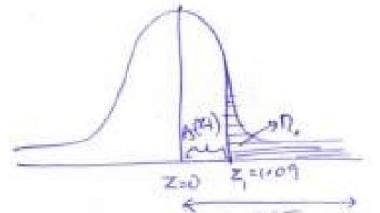
$$\Rightarrow A(z_1) + A(z < z_1) = 0.15$$

$$\Rightarrow A(z_1) = 0.5 - 0.1 = 0.4$$

From tables $z_1 = 1.28 \Rightarrow z_1 = -1.28$ { "left side of z" }

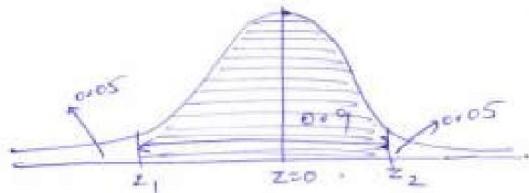
$$-1.28 = \frac{x-\mu}{\sigma} = \frac{x-0.78}{0.11} \Rightarrow x = 0.6392$$

Hence the highest mark obtained by the lowest 10% of students = $0.6392 \times 100\% = 64\%$.



Next page

P(iii)



To find $z_1 = ?$ $z_2 = ?$ Such that $A(z_1) + A(z_2) = 0.9$

Middle 90%. \Rightarrow Area in the middle $= 0.9$.

\Rightarrow Area on both sides left over $= 0.05$.

From table ! value of Z whose area $= 0.45 \Rightarrow Z = 1.64$

\therefore Middle area $= 0.9 \Rightarrow A(z_1) = A(z_2) = \frac{0.9}{2} = 0.45$

From table Area $= 0.45$ for $Z = 1.64$.

$\therefore z_1 = -1.64 ; z_2 = 1.64$.

$\Rightarrow \frac{z_1 - \mu}{\sigma} = -1.64 \Rightarrow \frac{z_2 - \mu}{\sigma} = 1.64$.

$\Rightarrow z_1 = \mu - 1.64\sigma ; z_2 = \mu + 1.64\sigma$.

$\therefore z_1 = \mu - 1.64\sigma ; \sigma = 0.11$.

Given $\mu = 78 ; \sigma = 0.9604$.

$\therefore z_1 = 0.5996 ; z_2 = 0.9604$.

\therefore Limits of 90% are $0.5996 \times 100 = 59.96 = 60\%$

& $0.9604 \times 100 = 96.04 = 96\%$.

\Rightarrow B/w 60 & 96 student marks.

8. If X is a normal variate with mean 30 & standard deviation 5, find the probabilities that (i) $26 \leq X \leq 40$ (ii) $X \geq 45$.

Sol: $\mu = 30, \sigma = 5 ; z = \frac{x-\mu}{\sigma}$

(i) $P(26 \leq X \leq 40) = ?$

$$z_1 = 26 \rightarrow z_1 = \frac{26-30}{5} = -\frac{4}{5} = -0.8$$

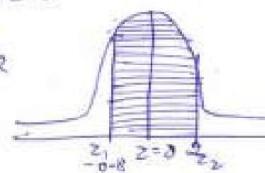
$$z_2 = 40 \rightarrow z_2 = \frac{40-30}{5} = \frac{10}{5} = 2$$

Thus $P(26 \leq X \leq 40) = P(-0.8 \leq z \leq 2)$

$$= A(0.8) + A(2)$$

$$= 0.2881 + 0.4772$$

$$= 0.7653 \text{ (using Normal distribution table)}$$



Problems on Joint Probability [ie Multiple Random Variables]

(1) Pg 237 Pb 2. For following bivariate p.d. of $X \& Y$ find

- (i) $P(X \leq 2, Y=2)$
- (ii) $P(X \leq 3, Y \leq 4)$
- (iii) $P(Y=3)$
- (iv) $F_X(2)$
- (v) $F_Y(3)$.

X/Y	1	2	3	4
1	0.1	0	0.2	0.1
2	0.05	0.12	0.08	0.01
3	0.1	0.05	0.1	0.09

$$(i) P(X \leq 2, Y=2) = P(X=1, Y=2) + P(X=2, Y=2)$$

$$= 0 + 0.12$$

$$= 0.12$$

$$(ii) P(X \leq 3, Y \leq 4) = P(X=1, Y=1) + P(X=1, Y=2) + P(X=1, Y=3)$$

$$+ P(X=1, Y=4) + P(X=2, Y=1) + P(X=2, Y=2)$$

$$+ P(X=2, Y=3) + P(X=2, Y=4)$$

$$= 0.1 + 0 + 0.2 + 0.1 + 0.05 + 0.12 + 0.08 + 0.01$$

$$= 0.66$$

$$(iii) P(Y=3) = P(X=1, Y=3) + P(X=2, Y=3) + P(X=3, Y=3)$$

$$= 0.2 + 0.08 + 0.1$$

$$= 0.38.$$

$$(iv) F_X(2) = P(X \leq 2) = P(X=1, Y=1) + P(X=1, Y=2) + P(X=1, Y=3) + P(X=2, Y=1) + P(X=2, Y=2) + P(X=2, Y=3) + P(X=2, Y=4)$$

$$= 0.05 + 0.12 + 0.08 + 0.01 + 0.1 + 0.2 + 0.1$$

$$= 0.66$$

$$(v) F_Y(3) = P(Y \leq 3) = P(X=1, Y=1) + P(X=1, Y=2) + P(X=1, Y=3) + P(X=2, Y=1) + P(X=2, Y=2) + P(X=2, Y=3) + P(X=3, Y=1) + P(X=3, Y=2) + P(X=3, Y=3)$$

$$= 0.1 + 0 + 0.2 + 0.05 + 0.12 + 0.08 + 0.1 + 0.05 + 0.1$$

$$= 0.70$$

INTRODUCTION

In a bivariate distribution and multivariate distributions we may be interested to find if there is any relationship between the two variables under study.

The Correlation is a statistical tool which studies the relationship between two variables and Correlation analysis involves various methods and techniques used for studying and measuring the extent of the relationship between them.

→ Two variables are said to be Correlated if the change in one variable results in a corresponding change in the other variable.

Types of Correlation1) positive and Negative Correlation

If the values of the two variables deviate in the same direction i.e. If the increase in the values of one variable results, on an average, in a corresponding increase in the value of the other variable (or) If the decrease in the value of one variable results, on an average, in a corresponding decrease in the value of the other variable, Correlation is said to be positive Correlation.

- Eg:-
1. Height and weights of the individuals.
 2. The family income and expenditure on luxury items
 3. Amount of Rainfall and yield of the Crop.

→ If the increase (decrease) in one Variable results, on an average, in a corresponding decrease (increase) in the value of the other variable, Correlation is said to be Negative Correlation

- Eg:-
1. Price and demand of a Commodity.
 2. Sale of winter garments and the day temperature
 3. Volume and pressure of a Perfect gas.

2. Linear and Nonlinear Correlation

→ The Correlation between two Variables is said to be linear if Corresponding to a unit change in one Variable, there is a constant change in the other Variable over the entire range of the values. (or)

Two Variables x and y are said to be linearly related, if there exists a relationship of the form $y = a + bx$.

→ Two Variables are said to be nonlinear or curvilinear if Corresponding to a unit change in one Variable, the other Variable does not change at a constant rate but at fluctuating rate.

(2)

Methods of studying Correlation

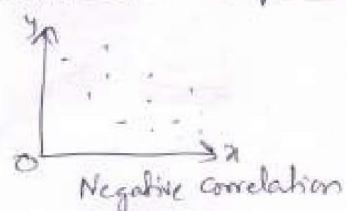
- 1) Scatter diagram method.
- 2) Karl Pearson's coefficient of correlation.
- 3) Rank method.

Scatter Diagram Method

If 'n' pair of values for two variables x and y are given as $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. Then these points are dotted on the x -axis and y -axis in the xy plane. The diagram of dots so obtained is known as scatter diagram.

Note:- This method gives a fairly good, though rough, idea about the relationship between the two variables.

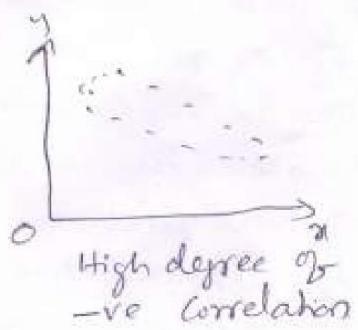
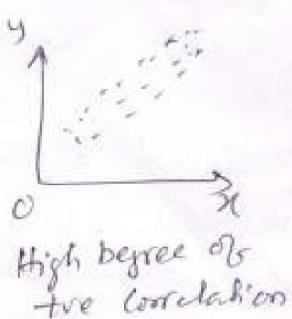
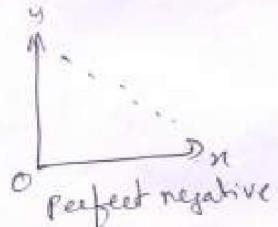
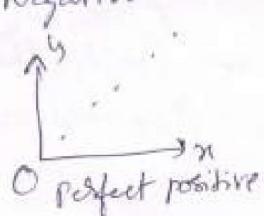
- If the points are very close (close to each other), then we say a fairly good amount of correlation between two variables.
- If the points are widely scattered, a poor correlation between those variables.
- If there is an upward trend rising from left lower left hand corner and going upward to the upper right hand corner, the correlation is positive.



→ If the points depict a downward trend from the upper left hand corner to the lower right hand corner, the correlation is negative.

→ If all the points, lie on straight line from the left bottom and getting going up towards the right top, the Correlation is said to be perfect positive

→ If all the points, lie on straight line from the left top and coming down toward right bottom, the correlation is said to be perfect negative



(3)

Prob1 Draw a scatter diagram from the following data

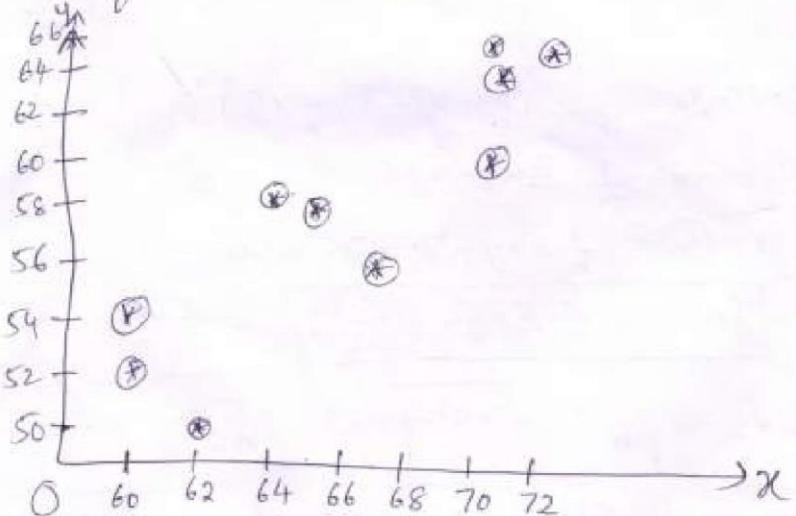
Height (inches) : 62 72 70 60 67 70 64 65 60 70

Weight (lbs) : 50 65 63 52 56 60 59 58 54 65

Also indicate whether Correlation is positive or negative.

Sol: Let us take Heights as random variable 'x' and taking those values on X-axis.

Let us take weights as variable 'y' and taking its values on Y-axis then we get scatter diagram as follows:



Since the points are close to each other we can say high degree of correlation (positive) between the series of heights and weights.

Karl Pearson's Coefficient of Correlation

Pearson's coefficient of correlation between two variables (series) X and Y is denoted by ' r ' is a measure of linear relationship between them and is defined as

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \quad \rightarrow ①$$

where $\text{Cov}(x, y) = \frac{1}{n} \sum (x - \bar{x})(y - \bar{y})$

$$\sigma_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} \text{ and } \sigma_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n}}$$

Note 1:

$$\text{Cov}(x, y) = \frac{1}{n} \sum (xy - \bar{x}\bar{y} - \bar{x}y + \bar{x}\bar{y})$$

$$\text{Cov}(xy) = \frac{1}{n} \sum xy - \frac{\sum x\bar{y}}{n} - \bar{x} \frac{\sum y}{n} + \bar{x}\bar{y}$$

$$\text{Cov}(x\bar{y}) = \frac{1}{n} \sum x\bar{y} - \bar{x}\bar{y} - \bar{x}\bar{y} + \bar{x}\bar{y}$$

$$\boxed{\text{Cov}(xy) = \frac{1}{n} \sum xy - \frac{\sum x(\bar{y})}{n}} \quad \rightarrow ②$$

Also $\sigma_x = \sqrt{\frac{1}{n} \sum (x^2 + \bar{x}^2 - 2x\bar{x})} = \sqrt{\frac{1}{n} \sum x^2 + \bar{x}^2 - 2\bar{x}^2}$

$$\therefore \sigma_x = \sqrt{\frac{1}{n} \sum x^2 - (\frac{\sum x}{n})^2} = \frac{1}{n} \sqrt{n \sum x^2 - (\sum x)^2}$$

$$\therefore \sigma_x = \frac{1}{n} \sqrt{n \sum x^2 - (\sum x)^2} \quad \rightarrow ③$$

$$\text{and } \sigma_y = \frac{1}{n} \sqrt{n \sum y^2 - (\sum y)^2} \rightarrow ④$$

Sub ②, ③, ④ in ① we get

$$\boxed{r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}} \quad \rightarrow ⑤$$

(4)

Note 2:-

$$r = \frac{\sum dx dy}{\sqrt{\sum dx^2 \sum dy^2}} \quad \text{where } dx = x - \bar{x} \quad \text{and} \quad dy = y - \bar{y}$$

- This formula is used when \bar{x} and \bar{y} are not fractions (i.e. not decimals or other natural numbers)
- If the values of x and y are large, \bar{x}, \bar{y} are fractions then using this formula becomes tedious. So, we use Step deviation method
- If the values of x and y are small but \bar{x}, \bar{y} are fractions then we use formula (5)

Properties of Correlation Coefficient

- * 1. The correlation coefficient lies between -1 and 1 i.e. $-1 \leq r \leq 1$.

Proof:- We know that $\sum \left[\frac{x-\bar{x}}{\sigma_x} \pm \frac{y-\bar{y}}{\sigma_y} \right]^2 \geq 0$

$$\sum \left\{ \left[\frac{x-\bar{x}}{\sigma_x} \right]^2 + \left[\frac{y-\bar{y}}{\sigma_y} \right]^2 \pm \frac{2(x-\bar{x})(y-\bar{y})}{\sigma_x \sigma_y} \right\} \geq 0$$

$$\Rightarrow \frac{\sum (x-\bar{x})^2}{\sigma_x^2} + \frac{\sum (y-\bar{y})^2}{\sigma_y^2} - \frac{2 \sum (x-\bar{x})(y-\bar{y})}{\sigma_x \sigma_y} \geq 0$$

Dividing with 'n' on both sides, we get

$$\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{n} + \frac{2 \text{cov}(x,y)}{\sigma_x \sigma_y} \geq 0 \quad [\because \frac{\sum (x-\bar{x})^2}{n} = \sigma_x^2]$$

$$2 + 2r \geq 0$$

$$\begin{aligned} 2 - 2r &\geq 0 \\ 2 &\geq 2r \Rightarrow 1 \geq r \\ \Rightarrow r &\leq 1 \end{aligned}$$

$$\begin{aligned} 2 + 2r &\geq 0 \\ 2 &\geq -2r \\ 1 &> -r \\ -1 &\leq r \end{aligned}$$

$$\therefore -1 \leq r \leq 1$$

\rightarrow If $r=1 \Rightarrow$ +ve correlation (perfect)
 $r=-1 \Rightarrow$ -ve correlation (perfect)
 $r=0 \Rightarrow$ Uncorrelated (No correlation)

2. The coefficient of Correlation is independent of the change of origin and scale.

Proof: let x and y are transformed into new variables u and v by the change of origin and scale

i.e
$$u = \frac{x-A}{h}, \quad v = \frac{y-B}{k}$$

where $h, k > 0$ and A, B, h, k are constants

then $\bar{x} = A + uh$ and $\bar{y} = B + vk$

$$\sum x = \sum (A + uh) \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{So we get}$$

$$\sum x = \frac{nA}{n} + \frac{h \sum u}{n} \quad \bar{y} = B + \bar{v}k$$

$$\Rightarrow \bar{x} = A + \bar{u}h$$

$$\text{Now } x - \bar{x} = h(u - \bar{u}) \text{ and } y - \bar{y} = k(v - \bar{v})$$

$$\therefore r_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} = \frac{\sum h k (u - \bar{u})(v - \bar{v})}{\sqrt{\sum (u - \bar{u})^2} \sqrt{\sum (v - \bar{v})^2}}$$

$$\therefore r_{xy} = \frac{\sum (u - \bar{u})(v - \bar{v})}{\sqrt{\sum (u - \bar{u})^2} \sqrt{\sum (v - \bar{v})^2}} = \frac{\text{cov}(uv)}{\sigma_u \sigma_v} = r_{uv}$$

$$\therefore r_{xy} = r_{uv}$$

Note:

$$r = \frac{\sum du dv}{\sqrt{\sum u^2} \sqrt{\sum v^2}} = \frac{n \sum uv - \sum u \sum v}{\sqrt{[n \sum u^2 - (\sum u)^2] [n \sum v^2 - (\sum v)^2]}} \quad \text{is}$$

used when \bar{x} or/and \bar{y} are fractionals or
if x, y values are large.

(5)

Problem 1. Calculate coefficient of correlation for the following data

x	9	8	7	6	5	4	3	2	1
y	15	16	14	13	11	12	10	8	9

Solution Given $n = 9$.

$$\therefore \bar{x} = \frac{\sum x}{n} = \frac{45}{9} = 5 \text{ and } \bar{y} = \frac{\sum y}{n} = \frac{108}{9} = 12.$$

Since the given values are small and \bar{x}, \bar{y} are not fractions,

correlation coefficient is given by

$$r = \frac{\sum dx dy}{\sqrt{\sum dx^2 \sum dy^2}} \quad \text{where } dx = x - \bar{x} \\ dy = y - \bar{y}$$

Now the table is

x	y	$dx = x - \bar{x}$	$dy = y - \bar{y}$	dx^2	dy^2	$dx dy$
9	15	4	3	16	9	12
8	16	3	4	9	16	12
7	14	2	2	4	4	4
6	13	1	-1	1	1	0
5	11	0	0	0	0	0
4	12	-1	0	4	4	4
3	10	-2	-2	9	16	12
2	8	-3	-4	9	16	12
1	9	-4	-3	16	9	12
$\sum x = 45$				$\sum dx^2 = 60$	$\sum dy^2 = 60$	$\sum dx dy = 57$

$$\therefore r = \frac{57}{\sqrt{60(60)}} = \frac{57}{60} = 0.95$$

\therefore There is a high degree of positive correlation between x and y series.

2. Find Correlation Coeff from

x	10	12	18	24	23	27
y	13	18	12	25	30	10

3. Find Karl Pearson's Coefficient of Correlation from

wages	100	101	102	102	100	99	97	98	96	95
Cost of living	98	99	99	97	95	92	95	94	90	91

Q. Given n = 10

$$\bar{x} = \frac{\sum x}{n} = \frac{2(100) + 101 + 2(102) + 99 + 97 + 98 + 96 + 95}{10} = \frac{990}{10} = 99$$

$$\bar{y} = \frac{\sum y}{n} = \frac{950}{10} = 95$$

The table is

x	y	$dx = \frac{x - \bar{x}}{x - 99}$	$dy = \frac{y - \bar{y}}{y - 95}$	dm	dy^2	$dx dy$
100	98	1	3	1	9	3
101	99	2	4	4	16	8
102	99	3	4	9	16	12
102	97	3	2	1	0	6
100	95	1	0	0	9	0
99	92	0	-3	4	0	0
97	95	-2	0	1	1	0
98	94	-1	-1	9	25	1
96	90	-3	-5	16	16	15
95	91	-4	-4	16	16	16
990	950				$\sum dx = 54$	$\sum dy = 96$
$\sum x$	$\sum y$					$\sum dx dy = 61$

$$\therefore r = \frac{\sum dx dy}{\sqrt{\sum dx^2 \sum dy^2}} = \frac{61}{\sqrt{54(96)}} = \frac{61}{\sqrt{5248}} = \frac{61}{72.04} = 0.847$$

$$\therefore r = 0.847$$

(4)

Q. Calculate the coeff of correlation between age of Cars and annual maintenance Cost and comment

Age of Cars (yr)	2	4	6	7	8	10	12
Annual Maintenance cost (Rs)	1600	1500	1800	1900	1700	2100	2000

SQ. Since 'y' values are large, we use step deviation method.

$$\text{Now, } n=7, \bar{x} = \frac{\sum x}{n} = \frac{49}{7} = 7$$

$$\text{and } \bar{y} = \frac{\sum y}{n} = \frac{12600}{7} = 1800$$

let $h=1$ and $K=100$ and $A=\bar{x}$; $B=\bar{y}$

$$\text{then } u = \frac{x-\bar{x}}{h} = x-7 \text{ & } v = \frac{y-1800}{100}$$

The table is:

x	y	u=x-7	v = $\frac{y-1800}{100}$	uv	u ²	v ²
2	1600	-5	-2	10	25	4
4	1500	-3	-3	9	9	9
6	1800	-1	0	0	1	0
7	1900	0	1	0	0	1
8	1700	1	-1	-1	1	1
10	2100	3	3	9	9	9
12	2000	5	2	10	25	4
49	12600	0	0	87	140	28
		Σu	Σv	Σuv	Σu^2	Σv^2

$$\therefore r = \frac{n \Sigma uv - \Sigma u \Sigma v}{\sqrt{[n \Sigma u^2 - (\Sigma u)^2][n \Sigma v^2 - (\Sigma v)^2]}} = \frac{7(37)}{\sqrt{7(70)(28)}} = \frac{37}{\sqrt{1960}} = 0.8357 \approx 0.836$$

Q. Find out the coeff of correlation form

Height of father (inches)	65	66	67	68	69	71	73
Height of son (inches)	67	68	64	72	70	69	70

Sol. Given $n=7$

$$\bar{x} = \frac{\sum x}{n} = \frac{65+66+67+68+69+71+73}{7} = 68.428$$

$$\bar{y} = \frac{\sum y}{n} = \frac{67+68+64+72+70+69+70}{7} = 68.57$$

Since \bar{x}, \bar{y} are fractions,

$$\text{let } u = \frac{x-68}{1} \text{ and } v = \frac{y-69}{1}. (\because h=l=1)$$

The table is

x	y	$u=x-68$	$v=y-69$	u^2	v^2	uv
65	67	-3	-2	9	4	6
66	68	-2	-1	4	1	2
67	64	-1	-5	1	25	5
68	72	0	3	0	9	0
69	70	1	1	1	1	1
71	69	3	0	9	0	0
73	70	5	1	25	1	5
		$\sum u = -3$	$\sum v = -3$	$\sum u^2 = 49$	$\sum v^2 = 41$	$\sum uv = 19$

$$\therefore r = \frac{n\sum uv - \sum u \sum v}{\sqrt{[n\sum u^2 - (\sum u)^2][n\sum v^2 - (\sum v)^2]}} = \frac{7(-19) + 9}{\sqrt{[7(49) - 9][7(41) - 9]}} \\ = \frac{142}{\sqrt{334(278)}} = \frac{142}{\sqrt{931852}} = \frac{142}{304.72} = 0.466$$

$$\therefore r = 0.466$$

(7)

7. Given $n=10$, $\sigma_x = 5.4$, $\sigma_y = 6.2$ and sum of the product of deviation from the mean of x and y is 66. find the correlation coeff.

Soln Given $n=10$, $\sigma_x = 5.4$, $\sigma_y = 6.2$ and $\sum (x-\bar{x})(y-\bar{y}) = 66$

$$\therefore r = \frac{\sum (x-\bar{x})(y-\bar{y})}{n\sigma_x \sigma_y} = \frac{66}{10(5.4)(6.2)} = \frac{66}{(5.4)(6.2)} = 0.197$$

Rank Correlation Coefficient

Spearman's rank correlation is given by

$$\rho = 1 - \frac{6 \sum d^2}{n(n^2-1)}$$

where $d_i = \text{Difference of corresponding ranks of } x \text{ and } y$.

$n = \text{no of the terms in the series}$

Properties

1. $-1 \leq \rho \leq 1$
2. If $\rho=1$, there is complete agreement in the order of the ranks and they are in same direction.
3. If $\rho=-1$, there is complete disagreement in the order of the ranks and they are in opposite direction.

Problems when ranks are given

1. Following are the ranks obtained by 10 students in two subjects, Statistics and Maths. To what extent the knowledge of the students in two subjects is related?

Statistics	1	2	3	4	5	6	7	8	9	10
Maths	2	4	1	5	3	9	7	10	6	8

SQ. Given $n=10$; let $x = \text{Ranks in Statistics}$
 $y = \text{Ranks in Maths}$

The table is

x	y	$d_i = x - y$	d_i^2
1	2	-1	1
2	4	-2	4
3	1	2	4
4	5	-1	1
5	3	2	4
6	9	-3	9
7	7	0	0
8	10	-2	4
9	6	3	9
10	8	2	4
		$\sum d_i^2 = 40$	

Rank Correlation

Coefficient is given by

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

$$= 1 - \frac{6(40)}{10(10^2-1)}$$

$$= 1 - \frac{24}{99}$$

$$= \frac{75}{99} = 0.7575$$

$$\therefore \rho = 0.76$$

2. A random sample of 5 students are selected and their grades in Maths and statistics are found to be

	1	2	3	4	5
Maths	85	60	73	40	90
Stats	93	75	65	50	80

(8)

S1: let x = Ranks in Maths.
 y = Ranks in Stats. $n = 5$

[Here give ~~1~~ ranks based on their marks in descending order]

Mark in Maths	Rank x	Mark in Stats	Rank y	$d_i = x - y$	d_i^2
85	2	93	1	1	1
60	4	75	3	1	1
73	3	65	4	-1	1
40	5	50	5	0	0
90	1	80	2	-1	1
					4

$$\therefore P = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} = 1 - \frac{6(4)}{5(5^2-1)} = 1 - \frac{6(4)}{5(24)}$$

$$\therefore P = 1 - \frac{1}{5} = \frac{4}{5} = 0.8$$

$$\boxed{\therefore P = 0.8}$$

Problems on equal or Repeated Ranks

formula

$$P = 1 - \frac{6}{n(n^2-1)} \left\{ \sum d_i^2 + \frac{1}{12} m(m^2-1) + \frac{1}{12} m(m^2-1) + \dots \right\}$$

Here $m \rightarrow$ no of times rank is repeated.

Procedure

→ If the items are repeated in the series then common ranks are assigned such that the mean of the ranks which these items would have got if they were different from each other and the next item will get the rank next to the rank used in computing the common rank.

Eg: 7. An item is repeated at rank 4, twice then

$$\text{Common Rank} = \frac{4+5}{2} = \frac{9}{2} = 4.5$$

and next rank will be '6'.

8. An Item is repeated thrice at rank 3 then

$$\text{Common Rank} = \frac{3+4+5}{3} = \frac{12}{3} = 4$$

and next rank will be '6'.

Prob1. The ranks of the 15 students in two subjects A and B are given below. The two numbers within the brackets denoting the ranks of the same student in A and B respectively.
 (1,10), (2,7), (3,2), (4,6), (5,4), (6,8), (7,3), (8,1), (9,11), (10,15), (11,7),
 (12,5), (13,14), (14,12) and (15,13).

Use Spearman's formula to find the rank correlation coefficient.

SQ: " Since ranks not repeated use

$$P = 1 - \frac{6 \sum d^2}{n(n^2-1)} \quad \underline{\text{Ans: } -0.514}$$

Where $n = 15$,

(3)

2. The following table gives the score obtained by 11 Students in English and Telugu. Find the rank Correlation Coeff.

Scores of English	40	46	54	60	70	80	82	85	85	90	95
Scores of Telugu	45	45	50	42	48	75	55	72	65	42	70

Qn. Given $n=11$

Marks in English	Scores in Telugu	Rank x	Rank y	d_i	d_i^2
40	45	11	7.5	3.5	12.25
46	45	10	7.5	2.5	6.25
54	50	9	6	3	9
60	43	8	11	-4	16
70	40	7	1	5	25
80	75	6	5	0	0
82	55	5	2	1.5	2.25
85	72	3.5	4	-0.4	0.16
85	65	3.5	10	-8	64
90	42	2	m ₂	-2	4
95	70	1	3		
				$\sum d_i^2$	134.91

∴ Spearman's rank Correlation is given by

$$\rho = 1 - \frac{6 \left\{ \sum d_i^2 + \frac{1}{12} m(m-1) + \frac{1}{12} m(m^2-1) \right\}}{n(n^2-1)}$$

$$= 1 - \frac{6 \left\{ 134.91 + \frac{1}{12} 2(3) \right\}}{11(120)}$$

$$= 1 - \frac{6 \left\{ 135.91 \right\}}{\frac{1820}{220}} = 1 - 0.618 = 0.382$$

$$\therefore \rho = 0.382$$

2. A sample of 12 fathers and their elder sons gave the following data about their elder sons. Calculate the Coeff of rank Correlation.

Fathers	65	63	67	64	68	62	70	66	68	67	69	71
Sons	68	66	68	65	69	66	68	65	71	67	68	70

SQP Here $n=12$. The table is

Father X	Sons Y	Rank x	Rank y	d_i	d_i^2
65	68	9	5.5	3.5	12.25
63	66	10	9.5	1.5	2.25
67	68	6.5	5.5	1	1
64	65	10	11.5	-1.5	2.25
68	69	4.5	3	1.5	2.25
62	66	12	9.5	2.5	6.25
70	68	2	5.5	-3.5	12.25
66	65	8	11.5	-3.5	12.25
68	71	4.5	1	3.5	12.25
67	67	6.5	9	-1.5	2.25
69	68	3	5.5	-2.5	6.25
71	70	1	2	-1	1
					72.5

In X series 68 repeated 2 times, 67 repeated 2 times

$$\text{So, } m = 2, 2.$$

In Y series 68 repeated 4 times, 65 repeated 2 times,
66 repeated 2 times

$$\text{So, } m = 4, 2, 2$$

$$\therefore \rho = 1 - \left(1 - \frac{\sum d + \frac{1}{12} \cdot 2(12-1)(2) + \frac{1}{12} \cdot 1(4^2-1)}{12(12-1)} \right) = 1 - \left(1 - \frac{72.5 + 2 + 5}{12(143)} \right)$$

$$\therefore \rho = 1 - \left(1 - \frac{79.5}{1716} \right) = 1 - \frac{79.5}{288} = 1 - 0.278 = \underline{\underline{0.722}}$$

$\therefore \rho = 0.722$

Regression

The statistical method which helps us to estimate the unknown value of one variable from the known value of the related variable is called Regression.

The lines described in the average relationship between two variables is known as lines of Regression.

Comparison between Correlation and Regression

<u>Correlation</u>	<u>Regression</u>
1. It is a measure of degree of covariability between two variables.	1. Regression establishes the functional relationship between dependent & independent Variable
2. In correlation, both the variables are random variables	2. In Regression, one variable is dependent Variable and other one is independent Variable
3. The coefficient of correlation is a relative measure	3. Regression Coefficient is an absolute measure

Linear Regression

→ If the Regression equation is a straight line then we say it is a Linear Regression.

Otherwise we call it as Non-linear Regression

Lines of Regression

- Line of Regression of y on x ^{is the line} which gives the best estimate for the value of y for any specified value of x .
- Line of Regression of x on y is the line which gives the best estimate for the value of x for any specified value of y .

Regression Equation of y on x $\rightarrow y = a + bx$ $\rightarrow \textcircled{1}$

$$\begin{aligned}na + b\sum x &= \sum y \\a\sum x + b\sum x^2 &= \sum xy\end{aligned}\quad (\text{Normal equations})$$

Solving we get a, b values

Substituting in $\textcircled{1}$ we get required
Regression eqn of y on x .

Regression Equation of x on y is $x = a + by$ $\rightarrow \textcircled{2}$

Normal equations are

$$\begin{aligned}na + b\sum y &= \sum x \\a\sum y + b\sum y^2 &= \sum xy\end{aligned}$$

Solving we get a, b values

Sub in $\textcircled{2}$ to get Required

Regression Equation of x on y .

(11)

→ Regression Equations by using deviations from arithmetic mean of x and y .

Regression Equation of ' y ' on ' x ' is

$$y - \bar{y} = r \frac{\sigma_x}{\sigma_y} (x - \bar{x})$$

where \bar{x} = Mean of ' x ' Series = $\frac{\sum x}{n}$

\bar{y} = Mean of ' y ' Series = $\frac{\sum y}{n}$

r = Correlation coefficient of x & y .

$\sigma_x, \sigma_y \rightarrow S.D. \text{ of } x, y \text{ series}$

Regression Coefficient of y on x = $r \frac{\sigma_y}{\sigma_x}$

$$= \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \cdot \frac{\sigma_y}{\sigma_x}$$

$$= \frac{\sum (x - \bar{x})(y - \bar{y})}{n \sigma_x^2}$$

$$= \frac{\sum xy}{n \sqrt{\sum (x - \bar{x})^2}}$$

$$\Rightarrow \text{Regression Coeff of } y \text{ on } x = \frac{\sum xy}{\sum x^2}$$

∴ Regression equation of ' y ' on x is

$$y - \bar{y} = \frac{\sum xy}{\sum x^2} (x - \bar{x})$$

∴ Regression Equation of ' x ' on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$\Rightarrow x - \bar{x} = \frac{\sum xy}{\sum y^2} (y - \bar{y})$$

Here if regression coeff of y on x is b_{yx} and
regression coeff of x on y is b_{xy} then

$$b_{yx} = \gamma \frac{\sigma_y}{\sigma_x} \text{ and } b_{xy} = \gamma \frac{\sigma_x}{\sigma_y}$$

$$\Rightarrow b_{yx} \cdot b_{xy} = \gamma^2$$

$$\Rightarrow \boxed{\gamma^2 = b_{xy} \cdot b_{yx}}$$

Note 1 Regression line always passes through (\bar{x}, \bar{y})

Prob If two regression lines of y on x and x on y
are $a_1x + b_1y + c_1 = 0$ and $a_2x + b_2y + c_2 = 0$ then
Prove that $a_1b_2 < a_2b_1$

Sol. Given lines can also be written as

Regression eqn of y on x is

$$y = -\frac{a_1}{b_1}x - \frac{c_1}{b_1} = -\frac{c_1}{b_1} + \left(-\frac{a_1}{b_1}\right)x$$

$$\Rightarrow y = a + bx$$

$$\text{Regression coeff of } y \text{ on } x = b = \left(-\frac{a_1}{b_1}\right) = b_{yx}$$

Similarly we can find.

$$\text{Regression coeff of } x \text{ on } y = \left(-\frac{b_2}{a_2}\right) = b_{xy}$$

$$\text{But } \gamma^2 = b_{yx} \cdot b_{xy} = \left(-\frac{a_1}{b_1}\right) \left(-\frac{b_2}{a_2}\right) = \frac{a_1 b_2}{a_2 b_1}$$

Since correlation coeff $-1 \leq \gamma \leq 1$

$$\Rightarrow \frac{a_1 b_2}{a_2 b_1} < 1 \Rightarrow \boxed{a_1 b_2 < a_2 b_1}$$

Angle between two regression lines

we know that

Regression equation of y on x is

$$y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$\Rightarrow \text{Slope of line } = m_1 = r \frac{\sigma_y}{\sigma_x}$$

Regression equation of x on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$\Rightarrow \text{Slope of line } = m_2 = \pm r \frac{\sigma_x}{\sigma_y}$$

$$\therefore \text{Angle between two lines} = \tan \theta = \frac{m_1 - m_2}{1 + m_1 m_2}$$

$$\therefore \tan \theta = \frac{r \frac{\sigma_y}{\sigma_x} - \pm r \frac{\sigma_x}{\sigma_y}}{1 + \frac{\sigma_y^2}{\sigma_x^2}}, \quad \text{or} \left[\tan^2 \theta \right]$$

$$= \frac{\frac{\sigma_y}{\sigma_x} (r - \pm \frac{1}{r})}{\left(\frac{\sigma_x^2 + \sigma_y^2}{\sigma_x^2} \right)}$$

$$\boxed{\tan \theta = \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \left(\frac{r^2 - 1}{r} \right)}$$

Note 1: If θ is acute, $\tan \theta = \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \left(\frac{1 - r^2}{r} \right)$

2. If θ is obtuse, $\tan \theta = \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \left(\frac{r^2 - 1}{r} \right)$.

3. If $r = 0 \Rightarrow \tan \theta = \infty \Rightarrow \boxed{\theta = \pi/2}$

\Rightarrow No relation between two regression variables

4. If $r = \pm 1 \Rightarrow$ the variables $\theta = 0$ or $\theta = \pi$ \Rightarrow Two lines are coincident or parallel

Pbl. Find the Regression line of x on y given

x	78	77	85	88	87	82	81	77	76	83	97	93
y	84	82	82	85	89	90	88	92	83	89	98	97

Soln. Regression equation of x on y is

$$x - \bar{x} = \frac{\Sigma xy}{\Sigma y^2} (y - \bar{y})$$

$$\Rightarrow x - \bar{x} = \frac{\Sigma xy}{\Sigma y^2} (y - \bar{y})$$

$$\text{Hence } \bar{x} = \frac{\Sigma x}{n} = \frac{78+77+\dots+93}{12} = \frac{1004}{12} = 83.66$$

$$\bar{y} = \frac{\Sigma y}{n} = \frac{84+82+\dots+97}{12} = \frac{1061}{12} = 88.41$$

Since \bar{x}, \bar{y} are not whole numbers, we use
Assumed mean method.

∴ Regression eqn of x on y is

$$x - \bar{x} = \frac{n \Sigma xy - \Sigma x \Sigma y}{n \Sigma y^2 - (\Sigma y)^2} (y - \bar{y})$$

x	y	$x = x - \bar{x}$ $= x - 84$	$y = y - \bar{y}$ $= y - 88.41$	y^2	xy
78	84	-6	-4	16	24
77	82	-7	-6	36	42
85	82	1	-6	36	-6
88	85	4	-3	9	-12
87	89	3	1	1	3
82	90	-2	2	4	-4
81	88	-3	0	0	0
77	92	-8	4	16	-28
76	83	-9	-5	25	40
83	89	-1	1	1	-1
97	98	13	10	100	130
93	99	9	11	121	99
1004	1061	-4	5	365	287

$$\therefore \text{Regression eqn of } x \text{ on } y \text{ is} \\ (x - 83.66) = \frac{12(287) + 20}{12(365) - 25} (y - 88.41) \Rightarrow \boxed{x = 0.795y + 12.38}$$

Problems on Regression

1. Find the most likely production corresponding to a rainfall 40 from the following data.

	Rainfall	Production
Average	30	500 kgs
S.D	5	100 kgs
r	0.8	

Soln we have to find Y value at $x = 40$.

so we have to find Regression eqn of Y on X

Given $\bar{x} = 30$, $s_x = 5$, $s_y = 100$; $\bar{y} = 500$; $r = 0.8$

Regression eqn of Y on X is

$$(Y - \bar{Y}) = r \frac{s_y}{s_x} (x - \bar{x})$$

$$\Rightarrow Y - 500 = 0.8 \left(\frac{100}{5} \right) (x - 30)$$

$$\Rightarrow Y - 500 = 8x - 240$$

$$\Rightarrow Y = 8x + 260$$

$$\therefore Y \text{ at } x = 40 \Rightarrow Y(40) = 8(40) + 260 \\ = 320 + 260$$

\therefore Production corresponding to a rainfall 40 = 480 kgs.

2. From a sample of 200 pairs of the observation the following quantities were calculated.

$$\sum x = 11.34, \sum y = 20.78, \sum x^2 = 12.16, \sum y^2 = 84.96, \sum xy = 22.13$$

then find Regression eqn of Y on X, compute the Correlation Coefficient.

Soln Regression eqn of Y on X is $y - \bar{y} = \frac{\sum xy}{\sum x^2} (x - \bar{x})$

$$\Rightarrow \left(Y - \frac{20.78}{200} \right) = \frac{22.13}{12.16} \left(x - \frac{11.34}{200} \right)$$

Problems on Regression

1. Find the most likely production corresponding to a rainfall 40 from the following data.

	Rainfall	Production
Average	30	500kg
S.D	5	100kg
r	0.8	

So we have to find Y value at $X = 40$.

So we have to find Regression eqn of Y on X

Given $\bar{X} = 30$, $s_x = 5$, $s_y = 100$; $\bar{Y} = 500$; $r = 0.8$

Regression eqn of Y on X is

$$(Y - \bar{Y}) = r \frac{s_y}{s_x} (X - \bar{X})$$

$$\Rightarrow Y - 500 = 0.8 \left(\frac{100}{5} \right) (X - 30)$$

$$\Rightarrow Y - 500 = 8(X - 30)$$

$$\Rightarrow Y = 8X + 260$$

$$\therefore Y \text{ at } X = 40 \Rightarrow Y(40) = 8(40) + 260 \\ = 320 + 260$$

. Production corresponding to a rainfall 40 = 480 kg.

2. From a sample of 200 pairs of the observation the following quantities were calculated.

$$\sum X = 11.34, \sum Y = 20.78, \sum X^2 = 12.16, \sum Y^2 = 84.96, \sum XY = 22.13$$

then find Regression eqn of Y on X, compute the Correlation Coefficient.

So Regression eqn of Y on X is $Y - \bar{Y} = \frac{\sum XY}{\sum X^2} (X - \bar{X})$

$$\Rightarrow \left(Y - \frac{20.78}{200} \right) = \frac{22.13}{12.16} \left(X - \frac{11.34}{200} \right)$$

The 8th performance for which Judge Q could not attend, was awarded 37 marks by judge P. If judge Q present, how many marks he will give to 8th perf? (15)
Ques. Let 'x' be the series of Marks by P and

'y' be the Series of Marks by Q.

$$\text{Here } \bar{x} = \frac{\sum x}{n} = \frac{46+42+44+40+43+41+45}{7} = 43.$$

$$\text{Also } \bar{y} = \frac{\sum y}{n} = 38.$$

Since \bar{x}, \bar{y} are whole numbers, we use the

Regression eqn of y on x as

$$y - \bar{y} = \frac{\sum xy}{\sum x^2} (x - \bar{x}) \quad \text{where } x = x - \bar{x} \\ y = y - \bar{y}$$

To find $\sum x^2, \sum xy$ the table is given as follows

x	$x = x - \bar{x}$ $= x - 43$	y	$y = y - \bar{y}$ $= y - 38$	x^2	xy
46	3	40	2	9	6
42	-1	38	0	1	0
44	1	36	-2	1	-2
40	-3	35	-3	9	9
43	0	35	-3	0	0
41	-2	37	-1	4	2
45	2	41	3	4	6
				$\sum x^2 = 28$	$\sum xy = 21$

∴ Regression Eqn of y on x is

$$y - 38 = \frac{21}{28} (x - 43)$$

$$\Rightarrow y = \frac{3x}{4} - \frac{129}{4} + 38$$

$$\Rightarrow y = 0.75x + 5.75$$

$$\text{For } x = 37 \Rightarrow y = 0.75(37) + 5.75 = 33.5$$

∴ If Q is present he will award 33.5 marks for the 8th performance.

4. The heights x of mothers & daughters are given in the following table. Estimate average height of daughter when the height of mother is 64.5 inches.

Height of mother (x)	62	63	64	64	65	66	68	70
Height of daughter (y)	64	65	61	69	67	68	71	65

Sol. Here we have to find y value at $x=64.5$, so we need to find Regression eqn of y on x

Now

$$\bar{x} = \frac{\sum x}{n} = \frac{62+63+64+64+65+66+68+70}{8} = 65.25$$

slly

$$\bar{y} = \frac{\sum y}{n} = 67$$

Since \bar{x} is not a whole number, we use Assume mean method to find Regression eqn of y on x .

\therefore Regression eqn of y on x is

$$y - \bar{y} = \frac{\sum y}{n} (x - \bar{x})$$

$$\Rightarrow y - \bar{y} = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sum x^2 - \frac{(\sum x)^2}{n}} (x - \bar{x})$$

$$\Rightarrow y - \bar{y} = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} (x - \bar{x})$$

x	y	$x = x - 65$	$y = y - 67$	x^2	xy
62	64	-3	-3	9	9
63	65	-2	-2	4	4
64	61	-1	-6	1	6
64	69	-1	2 ²⁰	1	-2
65	67	0	0	0	0
66	68	1	1	1	1
68	71	3	4	9	12
70	65	5	-2	25	-10
		$\sum x = 2$	$\sum y = -6$	$\sum x^2 = 50$	$\sum xy = 20$

\therefore Regression eqn of y on x is

(16)
0.48

$$(y - 67) = \frac{8(20) + 12}{8(50) - 4} (x - 65.25)$$

$$y = 67 + \frac{172}{396} (x - 65.25)$$

$$y = 0.4343x - 0.4344(65.25) + 67$$

$$\Rightarrow y = 0.4343x + 38.66$$

$$\text{For } x = 64.5 \Rightarrow y = 0.4343(64.5) + 38.66$$

$$y = 66.67$$

\therefore Height of daughter when mother height is
64.5 inches = 66.67 inches.

5. Determine the eqn of a straight line which fit the data

X	10	12	13	16	17	20	25
Y	10	22	24	27	29	33	37

Q6. Let the Regression line of y on x be $y = ax + bx$
then Normal eqns are $na + b\sum x = \sum y$
 $a\sum x + b\sum x^2 = \sum xy$

①

X	Y	x^2	xy
10	10	100	100
12	22	144	264
13	24	169	312
16	27	256	432
17	29	289	493
20	33	400	660
25	37	625	925
$\sum x = 113$		$\sum y = 184$	$\sum x^2 = 1938$
			$\sum xy = 3174$

Sub above values in ①

$$we \ get \ 7a + 113b = 182$$

$$113a + 1938b = 3174$$

Solving we get $a = 0.82$, $b = 1.56$

\therefore Regression eqn of y on x is

$$y = 0.82 + 1.56x$$

6. If θ is the angle b/w 2 regression lines and standard deviation of y is twice the S.D of x & $r = 0.25$ then find θ .

Sol: Given $\sigma_y = 2\sigma_x$ & $r = 0.25$ = Correlation Coeff of $x \& y$.

we know that

$$\begin{aligned}\tan \theta &= \left(\frac{1-r^2}{1+r^2} \right) \left(\frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \\ &= \left(\frac{1-0.25^2}{0.25} \right) \left(\frac{2\sigma_x^2}{5\sigma_x^2} \right) \\ &= 3.75 \left(\frac{2}{5} \right) = 1.5\end{aligned}$$

$$\therefore \theta = \tan^{-1}(1.5) = 56.36^\circ$$

$$\boxed{\therefore \theta = 56.36^\circ}$$

7. If $\sigma_x = \sigma_y = 5$, & the angle between regression lines is $\tan^{-1}(\frac{4}{3})$ then find r .

Sol: $\tan \theta = \left(\frac{1-r^2}{1+r^2} \right) \left(\frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right)$

$$\tan \theta = \left(\frac{1-r^2}{1+r^2} \right) \left(\frac{r^2}{25} \right)$$

$$\Rightarrow \tan \theta = \frac{1-r^2}{2r}$$

$$\Rightarrow \theta = \tan^{-1} \left(\frac{1-r^2}{2r} \right) = \tan^{-1} \left(\frac{4}{3} \right) \text{ (given)}$$

$$\Rightarrow \frac{1-r^2}{2r} = \frac{4}{3} \Rightarrow 3-3r^2 = 8r$$

$$\Rightarrow 3r^2 + 8r - 3 = 0$$

$$\Rightarrow r = \frac{-8 \pm \sqrt{64+36}}{6} = \frac{-8 \pm 10}{6} = \frac{-8+10}{6} \text{ or } \frac{2}{6}$$

$$\boxed{\therefore r = 0.33} \text{ Since } -1 < r < 1 \quad \therefore r \neq -3$$

(17)

8. If $x = 2y + 3$ & $y = kx + 6$ are the regression lines of x on y and y on x respectively then show that

$$a) 0 < k \leq \frac{1}{2}$$

$$b) \text{ If } k = \frac{1}{8} \text{ then find } r, \bar{x}, \bar{y}.$$

Sol Given $x = 2y + 3$, $y = kx + 6$

$$\Rightarrow r \frac{\partial x}{\partial y} = 2, \quad r \frac{\partial y}{\partial x} = k$$

$$\Rightarrow r = \frac{2 \partial y}{\partial x} \quad r = \frac{k \partial x}{\partial y}$$

$$\Rightarrow r^2 = 2k \quad (-1 \leq r \leq 1)$$

$$\Rightarrow 0 \leq 2k \leq 1 \quad \Rightarrow 0 \leq r^2 \leq 1$$

$$\Rightarrow 0 \leq k \leq \frac{1}{2}$$

$$\text{If } k = \frac{1}{8} \Rightarrow r^2 = 2 \left(\frac{1}{8} \right) \Rightarrow r = \pm \frac{1}{2}$$

Since (\bar{x}, \bar{y}) satisfies the regression eqns, we have

$$\bar{x} = 2\bar{y} + 3 \quad \text{&} \quad \bar{y} = \frac{1}{8}\bar{x} + 6$$

$$\bar{x} - 2\bar{y} = 3$$

$$\underline{- \bar{x} - 8\bar{y}} = \underline{-48}$$

$$6\bar{y} = 51$$

$$\Rightarrow \bar{y} = \frac{51}{6} \text{ sub in (1), we } \bar{x} = \frac{51}{3} + 3 = \frac{60}{3} = 20$$

$$\therefore \bar{x} = 20; \bar{y} = \frac{51}{6}$$

9. Calculate the coeff of correlation if the eqns of regression lines x on y and y on x are $7x - 16y + 9 = 0$ and $5y - 4x - 2 = 0$ and mean of n & y

Sol. Given Regression eqn of x on y & y on x , are

$$7x - 16y + 9 = 0 \Rightarrow x = \frac{16}{7}y + \frac{9}{7} \Rightarrow b_{xy} = \frac{16}{7}$$

$$-4x + 5y - 3 = 0 \Rightarrow y = \frac{4}{5}x + \frac{3}{5} \Rightarrow b_{yx} = \frac{4}{5}$$

Solving above eqns we get (\bar{x}, \bar{y}) since it satisfies the Regression lines.

So solving them we get

$$\bar{x} = \text{Solving them we get } \bar{x} = 0.1034, \bar{y} = 0.5172$$

$$\bar{y} =$$

$$\text{Correlation coeff } r = \sqrt{b_{xy} \cdot b_{yx}}$$

$$= \sqrt{\frac{16}{7} \times \frac{4}{5}} = 0.22$$

10. Given $N = 10$, $\sigma_x = 5.4$, $\sigma_y = 6.2$ & sum of the deviations product of
of ~~the~~ from the mean of x & y is 66. then find correlation coeff.

Sol? we know that $r = \sqrt{b_{xy} \cdot b_{yx}}$

$$\text{Where } b_{yx} = r \frac{\sigma_y}{\sigma_x} = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sigma_x^2} = \frac{66}{10(5.4)^2} = \frac{66}{10(5.4)^2}$$

$$\therefore b_{yx} = \frac{66}{10(5.4)^2} = 0.226$$

$$\text{So } b_{xy} = r \frac{\sigma_x}{\sigma_y} = \frac{\sum xy}{\sigma_y^2} = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sigma_y^2} = \frac{66}{10(6.2)^2} = \frac{66}{10(6.2)^2}$$

$$\Rightarrow b_{xy} = \frac{66}{10(6.2)^2} = \frac{66}{384.4} = 0.1716.$$

$$\therefore r = \sqrt{(0.226)(0.1716)} = 0.1969$$

$$\boxed{\therefore r = 0.1969}$$

Unit -3 :

Sampling distributions & Statistical Inferences .

Sampling :

Population : Population is totality of statistical data forming a subject of investigation .

eg: population of heights of Indians etc .

Size of population is no. of observations in population & is denoted by N .

Sampling : The process of selection of a sample is called Sampling.

eg : to assess the quality of a bag of rice , we examine only a portion of it by taking a handful of it ~~from~~ from a bag & then decide to purchase it or not .

Thus in estimating the characteristics of the population , instead of enumerating entire population , only the individuals in the sample are examined . Then the sample characteristics are utilised to estimate the population .

To eliminate any possibility of bias in the sampling procedure choose a random sample where observations are made independently & at random .

Sample : It is a subset of population & size of the sample is denoted by ' n ' .

Types of Sampling:

1. Random Sampling (or Probability Sampling) : It is the process of drawing a sample from a population in such a way that each member of the population has an equal chance of being included in the sample .
eg : Selecting randomly 20 words from a dictionary .

Note : (a) No. of samples with replacement = N^n
(b) " " " without " = N^C_n .

2. Stratified Random Sampling: The population is first subdivided into several parts (or small groups) called strata according to some relevant characteristics. Then a sample is selected from each stratum at random.
3. Systematic Sampling (or Quasi-Random Sampling): In this, population is arranged in some order. Then from first ' k ' items (say) one unit is selected at random. This unit & every k th unit combined together form systematic sample. In this method only the first member is chosen at random.
4. Purposive (or Judgement) Sampling: In this method, the units of sample are chosen according to convenience & personal choice of the individual, who selects the sample. This method is suitable when sample is small. Here, the investigator must have a thorough knowledge of the population. It is always subject to some kind of bias. eg: To select 20 students from a class of 100 to analyse extra-curricular activities of the students, the investigator would select those, who according to him, would represent the class.
5. Sequential Sampling: It consists of a sequence of sample drawn one after another from the population depending on results of previous samples.
ie. If first sample leads to no clear decision, a second sample is drawn, & if required a third sample is drawn to arrive at a final decision to accept or reject the lot.

Classification of Samples:

- (1) Large sample: If $n \geq 30$ sample is said to be large sample.
- (2) Small sample: If $n < 30$ sample is small sample or exact sample.

Parameter : It is a statistical measure based on all units of a population . eg: Mean μ , Variance σ^2 etc .

Statistic : It is a statistical measure based on all units of a selected sample . eg: sample mean \bar{x} , sample variance s^2 etc .

Note : Value of a statistic varies from sample to sample (\because units of samples are not same), but parameters always remains constant .

SAMPLING DISTRIBUTION of a STATISTIC

It helps us

- (1) To estimate unknown population parameter from known statistic .
- (2) To set confidence limits of parameter within which the parameter values are expected to lie .
- (3) To test a hypothesis & to draw a statistical inference from it .

Central Limit theorem

If \bar{x} be mean of a sample size n , drawn from a population with mean μ & S.D σ then sampling dist of sample mean \bar{x} is approx normal dist with mean μ , $S.E = \frac{\sigma}{\sqrt{n}}$ (n ≥ 30)

$Z = \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)}$ is a random variable whose distribution function approaches that of standard normal distribution $N(0,1)$ as $n \rightarrow \infty$.

Standard Error (S.E) : It enables us to determine the confidence limits within which the parameters are expected to lie .

Imp. Formulae :

$$(1) S.E \text{ of sample mean} = \frac{\sigma}{\sqrt{n}}$$

$$(2) " " " \text{ proportion} = \sqrt{\frac{PQ}{n}}, \text{ where } Q = 1 - P$$

$$(3) " " " \text{ Standard deviation} = \frac{\sigma}{\sqrt{2n}}$$

(4) S.E of difference of 2 sample means \bar{x}_1 & \bar{x}_2 is = $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$
when populations are different.

(b) when samples are from same population then

$$S.E \text{ of difference of 2 sample means } \bar{x}_1, \bar{x}_2 = \sqrt{\sigma^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

(c) when samples are having standard deviations s_1, s_2 for the 2 samples whose means are \bar{x}_1, \bar{x}_2 then

$$S.E \text{ of 2 sample means } \bar{x}_1, \bar{x}_2 = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\text{or} = \sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \text{ where}$$

$$s^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2}.$$

(5) S.E of difference of 2 sample proportions p_1, p_2

$$= \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}} \text{ where } p_1, p_2 \text{ are population proportions}$$

$$= \sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \text{ when samples are from same population}$$

$$= \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}} \text{ when population proportions are not given.}$$

$$= \sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \text{ where } P = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \text{ &} \\ Q = 1 - P.$$

Correction factor = $\frac{N-n}{N-1}$

For finite population of size N , when a sample is drawn without replacement,

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

$$(ii) S.E of proportion = $\sqrt{\frac{pq}{n}} \cdot \sqrt{\frac{N-n}{N-1}}$.$$

Sampling Distribution of Mean (त्रिकुम्भ)

(i) Mean of the population: $\mu = \frac{\sum x_i}{N}$

(ii) Variance " " " $\sigma^2 = \frac{1}{N} \sum (x_i - \bar{x})^2$

Do problems.

- Ques: A population consists of 5, 10, 14, 18, 13, 24. Consider all possible samples of size two which can be drawn without replacement from the population. Find
- Mean of the population
 - S.D " "
 - mean of sampling distribution of means.
 - S.D " "

Sol: To make samples of size 2 without replacement:

(5, 10), (5, 14), (5, 18), (5, 13), (5, 24)
 (5, 10), (5, 14), (5, 18), (10, 13), (10, 24),
 (10, 14), (10, 18), (10, 13), (14, 13), (14, 24),
 (14, 18), (14, 13), (18, 24)
 (18, 13), (18, 24)

$$\text{No. of samples} = {}^6 C_2 (\because N=6) = \frac{6 \times 5}{2} = 15$$

(a) Mean of population $\mu = \frac{5+10+14+18+13+24}{6} \quad \{ \because N=6 \}$
 $= \frac{88}{6} = 14 \quad \therefore \mu = 14$

(b) S.D of population $\sigma = \sqrt{\sigma^2}$ where
 $\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N} = \frac{(5-14)^2 + (10-14)^2 + (14-14)^2 + (18-14)^2 + (13-14)^2 + (24-14)^2}{6}$
 $= \frac{(5-14)^2 + (10-14)^2 + (14-14)^2 + (18-14)^2 + (13-14)^2 + (24-14)^2}{6}$
 $= \frac{81+16+0+16+1+100}{6}$
 $\Rightarrow \sigma^2 = \frac{214}{6} = 35.67$
 $\therefore \sigma = \sqrt{35.67}$

Mean of the samples :

$$\frac{5+10}{2}, \frac{5+14}{2}, \frac{5+18}{2}, \frac{5+13}{2}, \frac{5+24}{2}$$

$$\frac{10+14}{2}, \frac{10+18}{2}, \frac{10+13}{2}, \frac{10+24}{2}$$

$$\frac{14+18}{2}, \frac{14+13}{2}, \frac{14+24}{2}$$

$$\frac{18+13}{2}, \frac{18+24}{2}$$

$$\frac{13+24}{2}$$

ie	7.5	9.5	11.5	9	14.5
	12	14	11.5	17	
	16	13.5	19		
	15.5	21			
	18.5				

c) Mean of sampling distribution of mean is $\mu_{\bar{x}}$

$$\mu_{\bar{x}} = \frac{(7.5+9.5+11.5+9+14.5+12+14+11.5+17+16+18.5+19+15.5+21)}{15} (18.5)$$

$$= 14.$$

d) Variance of S.D. of mean is

$$\sigma^2_{\bar{x}} = \frac{(7.5-14)^2 + (9.5-14)^2 + (11.5-14)^2 + (9-14)^2 + (14.5-14)^2 + (12-14)^2 + (16-14)^2 + (11.5-14)^2 + (17-14)^2 + (16-14)^2 + (13.5-14)^2 + (19-14)^2 + (21-14)^2 + (18.5-14)^2}{15}$$

$$= \frac{214}{15} = 14.2666$$

i. S.D. of sampling distribution of means

$$\sigma_{\bar{x}} = \sqrt{14.2666} = 3.78.$$

Ques. Samples of size 2 are taken from population 3, 6, 9, 15, 27 with replacement. Find (a) μ (b) σ^2 (c) $\mu_{\bar{x}}$ (d) $\sigma_{\bar{x}}^2$.

$$(a) \mu = \frac{3+6+9+15+27}{5} = 12. \quad (N=5)$$

$$(b) \sigma^2 = \frac{(3-12)^2 + (6-12)^2 + (9-12)^2 + (15-12)^2 + (27-12)^2}{5} = \frac{360}{5} = 72$$

$$\therefore \sigma = \sqrt{72}$$

(c) Samples of size 2 with replacement = $N^n = 5^2 = 25$. They are:
 $(3,3) (3,6) (3,9) (3,15) (3,27)$
 $(6,6) (6,3) (6,9) (6,15) (6,27)$
 $(9,9) (9,3) (9,6) (9,15) (9,27)$
 $(15,15) (15,3) (15,6) (15,9) (15,27)$
 $(27,27) (27,3) (27,6) (27,9) (27,15)$

Mean $\mu_{\bar{x}} = ?$

Means of samples are

3	4.5	6	9	15
6	4.5	7.5	10.5	16.5
9	6	7.5	12	18
15	9	10.5	12	21
27	15	16.5	18	21

$$\therefore \mu_{\bar{x}} = \frac{[3+4.5+6+9+15+6+4.5+7.5+10.5+12.5+9+6+7.5+\dots+(21-12)^2]}{25}$$

$$= \frac{300}{25} = 12.$$

$$(d) \sigma_{\bar{x}}^2 = \frac{(3-12)^2 + (4.5-12)^2 + (6-12)^2 + (9-12)^2 + (15-12)^2 + (6-12)^2 + (4.5-12)^2 + (10.5-12)^2 + (16.5-12)^2 + (7.5-12)^2 + (12-12)^2 + (18-12)^2 + (9-12)^2 + (10.5-12)^2 + (12-12)^2 + (21-12)^2 + (27-12)^2 + (15-12)^2 + (16.5-12)^2 + (18-12)^2 + (21-12)^2}{25}$$

$$= \frac{909}{25}$$

$$\therefore \sigma_{\bar{x}} = \sqrt{\frac{909}{25}} = 6.03.$$

- ③ A random sample of size 100 is taken from an infinite population having $\mu = 76$ & variance $\sigma^2 = 256$. What is the probability that \bar{x} will be b/w 75 & 78.

Sol: $n = 100$, $\mu = 76$, $\sigma^2 = 256 \Rightarrow \sigma = 16$.

To find $P(75 \leq \bar{x} \leq 78) = ?$

$$\text{we have } z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$\bar{x}_1 = 75 \Rightarrow z_1 = \frac{75 - 76}{16/\sqrt{100}} = -0.625$$

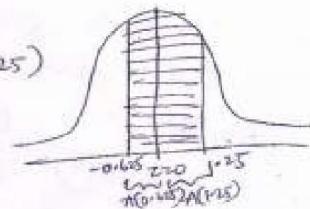
$$\bar{x}_2 = 78 \Rightarrow z_2 = \frac{78 - 76}{16/\sqrt{100}} = 1.25$$

$$\therefore P(75 \leq \bar{x} \leq 78) = P(-0.625 \leq z \leq 1.25)$$

$$= A(0.625) + A(1.25)$$

$$= 0.2324 + 0.3944$$

$$= 0.628 \text{ (ans.)}$$



- ④ A normal population has mean of 0.1 & S.D of 2.1. Find the probability that mean of a sample of size 900 will be negative.

Sol: $\mu = 0.1$, $\sigma = 2.1$, $n = 900$. To find $P(\bar{x} < 0) = ?$

$$\text{we have } z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{\bar{x} - 0.1}{2.1/\sqrt{900}} = \frac{\bar{x} - 0.1}{0.07}$$

$$\bar{x} - 0.1 = 0.07z \Rightarrow \bar{x} = 0.1 + 0.07z$$

$$\therefore P(\bar{x} < 0) = P(0.1 + 0.07z < 0) = P(0.07z < -0.1)$$

$$= P(z < -\frac{0.1}{0.07})$$

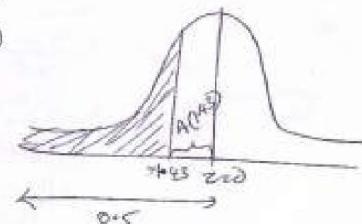
$$= P(z < -1.43)$$

$$\approx 0.5 - A(1.43)$$

$$\approx 0.5 - 0.4236$$

$$= 0.0764$$

Ans.



ESTIMATION

They are of 2 types: ① Point Estimation ② Interval Estimation.

① Defn: If an estimate of the population parameter is given by a single value, then the estimate is called a Point Estimation of the parameter. eg: If height of a student is measured as 162 cms then it gives a point estimation.

② If an estimate of the population is given by two different values b/w which the parameter may be considered to lie, then the estimate is called an Interval estimation of the parameter.

eg: If height of student lies b/w 159.5 cm & 166.5 cm then it gives interval estimation.

Confidence interval estimates of parameters

The confidence interval has probability of correctly estimating the true value of the population parameter.

$$\begin{aligned} \bar{x} &\pm 1.96 (\text{S.E. of } \bar{x}) \\ p &\pm 1.96 (\text{S.E. of } p) \\ (\bar{x}_1 - \bar{x}_2) &\pm 1.96 (\text{S.E. of } (\bar{x}_1 - \bar{x}_2)) \\ (p_1 - p_2) &\pm 1.96 (\text{S.E. of } (p_1 - p_2)) \end{aligned}$$

95% are confidence limits for
single mean
" proportion
Difference of mean
" " proportion,

Maximum Error of Estimate E :

$$\boxed{E = z_{\alpha/2} \left(\frac{\sigma}{\sqrt{n}} \right)}$$

for large sample

$$\boxed{E = t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right)}$$

" small n ~

Do problems:

Testing of Hypothesis

There are many problems, in which, rather than estimating the value of a parameter we need to decide whether to accept or reject a statement about the parameter.

This statement is called a hypothesis & the decision-making procedure about the hypothesis is called Testing of hypothesis. The procedure which enables us to decide on the basis of sample results whether a hypothesis is true or not, is called Test of Hypothesis or Test of Significance.

Test of Hypothesis involves following steps :

Step 1 : Statement (or assumption) of hypothesis ,

H_0 (Null hypothesis) &

H_1 (Alternative hypothesis) .

Step 2 : Specification of Level of Significance (α) :

L.O.S is the max. possibility with which we are willing to risk an error in rejecting H_0 when it is true.

Step 3 : Identification of Test Statistic :

say z, t, F etc .

Step 4 : Critical region : ie z_c or t_c etc .

Step 5 : Making Decision

If computed value $<$ critical value, we

Accept H_0 , otherwise reject H_0 ,

These are 2 types of hypothesis

(1) Null hypothesis : Denoted by H_0 . It is a definite statement which asserts that there is no significant difference b/w the statistic & the population parameter

eg : $H_0 : \mu = \mu_0$

(2) Alternative hypothesis : It always contradicts H_0 & is denoted by H_1 . If $H_0 : \mu = \mu_0$ then H_1 would be

(i) $H_1 : \mu \neq \mu_0$ (ii) $H_1 : \mu > \mu_0$ (iii) $H_1 : \mu < \mu_0$

(i) is known as 2-tailed alternative & (ii) is right tailed
(iii) is " " left - "

Errors of Sampling : They are

(i) Type I error : Reject H_0 when it is true

(or)

α error

(or)

Producer's risk

(ii) Type II error (or) β error (or) Consumer's risk :

Accept H_0 when it is wrong.

Critical Values (Z_d) of Z (i.e. for Large samples where $z = \frac{t - E(t)}{S.E.(t)}$)

	1% (0.01)	5% (0.05)	10% (0.1)
Two-tailed Test	$ Z_d = 2.58$	$ Z_d = 1.96$	$ Z_d = 1.645$
Right - " "	$Z_d = 2.33$	$Z_d = 1.645$	$Z_d = 1.28$
Left - " "	$= -2.33$	$= -1.645$	$= -1.28$

Test of Significance for Large Samples

Tests of significance used in large samples are different from those used in small samples, because if S.D. of population is not known, it can be replaced by S.D. of sample, whereas in small samples it is not possible.

Under Large samples we see 4 important tests of significance:

(1) Testing of significance for single mean, difference of means -

(2) " " " " " single proportion,

(3) " " " " " difference of proportions.

(4) " " " " " difference of proportions.

Note: The Confidence limits (or fiducial limits) for single mean

are $\bar{x} \pm 2.58 \frac{\sigma}{\sqrt{n}}$ or $(\bar{x} - 2.58 \frac{\sigma}{\sqrt{n}}, \bar{x} + 2.58 \frac{\sigma}{\sqrt{n}})$ for

99% confidence

(ii) $\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}}$ for 95% confidence limits; etc.

Formulae for Testing of Hypothesis

Test of Hypothesis for single mean:

$$(i) z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \quad \text{where } \bar{x} \rightarrow \text{sample mean}$$

$\mu \rightarrow \text{population's "}$

$\sigma \rightarrow \text{" S.D}$

$n \rightarrow \text{sample size.}$

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

(For) Testing of Hypothesis for difference of means:

$$(i) z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad \text{where } \bar{x}_1, \bar{x}_2 \rightarrow \text{sample means}$$

$\sigma_1, \sigma_2 \rightarrow \text{S.D of populations}$

$n_1, n_2 \rightarrow \text{sample sizes}$

(Here $S.E = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$)

(ii) If samples are from same population then

$$\sigma_1^2 = \sigma_2^2 = \sigma^2 \quad \& \quad z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\sigma^2(\frac{1}{n_1} + \frac{1}{n_2})}} \quad \text{(Here } S.E = \sqrt{\sigma^2(\frac{1}{n_1} + \frac{1}{n_2})} \text{)}$$

(iii) If S.D of populations are unknown they can be replaced by S.D s_1, s_2 of samples.

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

Note: The confidence limits are

a) $(\bar{x}_1 - \bar{x}_2) \pm 2.58 (S.E \text{ of } \bar{x}_1 - \bar{x}_2)$ for 99%.

b) $(\bar{x}_1 - \bar{x}_2) \pm 1.96 (S.E \text{ of } (\bar{x}_1 - \bar{x}_2))$ for 95%.

For Single Proportion:

$$z = \frac{p - P}{\sqrt{\frac{PQ}{n}}} \quad \text{where } P = \frac{x}{n} \rightarrow \text{sample proportion}$$

$P \rightarrow \text{Population's "}$

$Q = 1 - P$

$n \rightarrow \text{sample size.}$

$x \rightarrow \text{no. of successes.}$

Note: Confidence interval for proportion P are:

$$(P \pm 3\sqrt{\frac{PQ}{n}})$$

If P is not known we use ' p ' (sample proportion)

& limits are $(p \pm 3\sqrt{\frac{pq}{n}})$.

For difference of proportions:

$$(i) z = \frac{p_1 - p_2}{\sqrt{\frac{p_1 q_1 + p_2 q_2}{n_1 + n_2}}} \quad \text{where } p_1, p_2 \rightarrow \text{sample proportions}$$

$\underbrace{\qquad\qquad\qquad}_{S.E.}$

$p_1, p_2 \rightarrow \text{Population proportions}$
 $q_1 = 1 - p_1, q_2 = 1 - p_2$
 $n_1, n_2 \rightarrow \text{sample sizes}$
Also $p_1 = \frac{x_1}{n_1}, p_2 = \frac{x_2}{n_2}$

Note: When population proportions are unknown use sample proportions. i.e. $z = \frac{p_1 - p_2}{\sqrt{\frac{p_1 q_1 + p_2 q_2}{n_1 + n_2}}} \quad \underbrace{\qquad\qquad\qquad}_{S.E.}$

$$(ii) \text{Other method (called method of pooling)} \quad z = \frac{p_1 - p_2}{\sqrt{\frac{pq(L+\frac{1}{n_1+n_2})}{n_1+n_2}}} \quad \text{where } P = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \left(= \frac{x_1 + x_2}{n_1 + n_2}\right)$$

$$\text{or } z = \frac{p_1 - p_2}{\sqrt{\frac{pq(L+\frac{1}{n_1+n_2})}{n_1+n_2}}} \quad \text{where } p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \left(= \frac{x_1 + x_2}{n_1 + n_2}\right).$$

The confidence interval or limits are:

$$\textcircled{a} (p_1 - p_2) \pm z_{\alpha/2} \left[\text{S.E. of } p_1 - p_2 \right]$$

Do Problems.

① A truck company claims that average life of certain types is atleast 88000 miles. To check the claim it puts 40 of these tyres on its trucks & gets a mean life time of 27463 miles with a S.D of 1348 miles. Can the claim be true?

Sol: $n = 40, \bar{x} = 27463, s = 1348, \mu = 28000 \quad = \frac{1348}{\sqrt{40}} (\because \sigma \text{ is not given})$

$$H_0: \mu = 28000$$

$$H_1: \mu \neq 28000 \quad (\text{2 tailed test})$$

$$\alpha = 0.05 \quad (\because \text{not mentioned})$$

$$z_{\alpha/2} = 1.96 \quad (\text{for 95% confidence})$$

$$z = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{27463 - 28000}{1348/\sqrt{40}} = -2.52$$

$$\therefore |z| = 2.52 \quad \text{&} \quad |z| > z_{\alpha/2} \Rightarrow \text{Reject } H_0$$

i.e. The claim is not true.

Solved important model questions of Unit-III

1. A population consists of 5, 10, 14, 18, 13, 24. Consider all possible samples of size 2 from which can be drawn without replacement from the population. Find
 (a) Mean of the population (b) Standard deviation of population
 (c) Mean of Sampling distribution of means.
 (d) Mean of Sampling distribution of means.

Ex. Given population size = $N = 6$.

Sample size $n = 2$.

No. of samples of size '2' selected from given population
 (without replacement) = $N_{C_2} = 6C_2 = \frac{6!}{4!2!} = \frac{6 \times 5}{2 \times 1} = 15$

and The samples are.

$$\begin{cases} (5, 10), (5, 14), (5, 18), (5, 13), (5, 24) \\ (10, 14), (10, 18), (10, 13), (10, 24) \\ (14, 13), (14, 12), (14, 24) \\ (13, 12), (18, 24), (13, 24) \end{cases}$$

$$(a) \text{Mean of the population} = \mu = \frac{5+10+14+18+13+24}{6} = \frac{84}{6} = 14$$

$$\therefore \mu = 14$$

$$(b) \text{Standard Deviation of the population} = \sigma = \sqrt{\frac{\sum (x - \mu)^2}{n}}$$

$$\Rightarrow \sigma^2 = \frac{\sum (x - \mu)^2}{N} = \frac{(5-14)^2 + (10-14)^2 + (18-14)^2 + (13-14)^2 + (24-14)^2}{6}$$

$$\Rightarrow \sigma^2 = \frac{81 + 16 + 16 + 1 + 100}{6} = \frac{133 + 81}{6} = \frac{214}{6} = 35.67$$

$$\therefore \sigma = \sqrt{35.67} = 5.972$$

$$\therefore \sigma = 5.972$$

c) The means of samples of size 2 are

$$\{7.5, 9.5, 11.5, 9, 14.5, 12, 14, 11.5, 17, 16, 13.5, 19, 15.5, 21, 18.5\}$$

$$\therefore \text{The mean of sample means} = \frac{210}{15} = 14$$

$$\therefore \text{The Mean of the Sampling distribution of means} = 14 = \bar{\mu}_{\bar{x}}$$

$$\therefore \text{we observed that } \boxed{\mu = \mu_{\bar{x}}}$$

d) The Variance of sampling distribution of mean is

$$\sigma_{\bar{x}}^2 = \frac{(7.5-14)^2 + (9.5-14)^2 + \dots + (18.5-14)^2}{15}$$

$$= \frac{42.25 + 20.25 + 6.25 + \dots + 20.25}{15} = 14.2666$$

$$\therefore \sigma_{\bar{x}} = \sqrt{14.266} = 3.78$$

$$\therefore \text{S.D of sampling distribution of means} = 3.78$$

Technique

S.No	Sample	Sample mean \bar{x}	Mean of Sampling Distribution $\bar{\mu}_{\bar{x}} = \bar{x} - 14$	$(\bar{x} - \bar{\mu})^2$
1	(5, 10)	7.5	-6.5	42.25
2	(5, 14)	9.5	-4.5	20.25
3	(5, 18)	11.5	-2.5	6.25
4	(5, 12)	9	-5	25
5	(5, 24)	14.5	0.5	0.25
6	(10, 14)	12	-2	4
7	(10, 18)	14	0	0
8	(10, 13)	11.5	-2.5	6.25
9	(10, 4)	7	3	9
10	(14, 11)	16	2	4
11	(14, 13)	13.5	-0.5	0.25
12	(14, 24)	19	5	25
13	(18, 13)	15.5	1.5	2.25
14	(18, 24)	21	7	49
15	(13, 24)	18.5	4.5	20.25
Total		$\sum \bar{x} = 210$		214
			\downarrow	$\sum (\bar{x} - \bar{\mu})^2$

2. A population consists of five numbers 2, 3, 6, 8 and 11. Consider all possible samples of size two which can be drawn with replacement from this population. And
- The mean of the population
 - S.D of a population
 - Mean of the Sampling distribution of means.
 - Standard Error of Means. (S.D of the Sampling distribution of means)

Sol. Given $N=5$; (Population Size)
 \bar{x}) Mean of the population = $\frac{2+3+6+8+11}{5} = \frac{30}{5} = 6$

a) Mean of the population (σ^2) is

b) Variance of the population (σ^2) is

$$\sigma^2 = \frac{\sum (x-\mu)^2}{N} = \frac{1}{5} [(2-6)^2 + (3-6)^2 + (6-6)^2 + (8-6)^2 + (11-6)^2]$$

$$\therefore \sigma^2 = 10.8 \Rightarrow \sigma = \sqrt{10.8} = 3.29 \quad (3)$$

c) S.D of the population = 3.29.

c) Number of samples of size '2' from given population

$$N^n = 5^2 = 25$$

and the samples are given in the table

Sample	Mean of sample (\bar{x})	$(x - \bar{x})^2$
(2,2)	2	4
(2,3)	2.5	6.25
(2,6)	4	16
(2,8)	5	25
(2,11)	6.5	0.25
(3,2)	2.5	6.25
(3,8)	3	9
(3,6)	4.5	2.25
(3,8)	5.5	0.25
(3,11)	7	1
(6,2)	4	4
(6,3)	4.5	2.25
(6,6)	6	0
(6,8)	7	1
(6,11)	8.5	6.25

Sample	Mean of sample	$\frac{(x - \bar{x})^2}{n-1}$
(8, 2)	5	1
(8, 3)	5.5	0.25
(8, 6)	7	1
(8, 8)	8	4
(8, 11)	9.5	12.25
(11, 2)	6.5	0.25
(11, 3)	7	1
(11, 6)	8.5	6.25
(11, 8)	9.5	12.25
(11, 11)	11	25
Total	150	135

From the table, we have

(4)

$$\sum n = 150$$

$$\Rightarrow M_{\bar{x}} = \frac{\sum x}{n} = \frac{150}{25} = 6$$

$$\text{and } \sum (x - \bar{x})^2 = 135$$

$$\Rightarrow S_{\bar{x}}^2 = \frac{1}{n} \sum (x - \bar{x})^2 = \frac{135}{25} = 5.4$$

$$\therefore S_{\bar{x}} = \sqrt{5.4} = 2.32.$$

3. If the population is 3, 6, 9, 15, 27

a) List all possible samples of size 3 that can be taken without replacement from the finite population.

b) Mean of Sampling distribution of means

c) Find S.D of Sampling distribution of means

Sol. No. of samples of size 3 from given population = $N_{n,3} = {}^5C_3$

$$= \frac{5!}{2!3!} = \frac{5 \times 4}{2} = 10 \quad (\text{Without Replacement})$$

The samples are { (3, 6, 9), (3, 6, 15), (3, 6, 27), (3, 9, 15), (3, 9, 27) }

(3, 15, 27), (6, 9, 15), (6, 9, 27), (6, 15, 27)

(9, 15, 27) }

Sample	Mean of Sample	$\frac{(x - \bar{x})^2}{n-12}$
(3, 6, 9)	6	3.6
(3, 6, 15)	8	1.6
(3, 6, 27)	12	0
(3, 9, 15)	9	9
(3, 9, 27)	13	1
(3, 15, 27)	15	9
(6, 9, 15)	10	4
(6, 9, 27)	14	4
(6, 15, 27)	16	1.6
(9, 15, 27)	17	25
	$\Sigma x = 120$	$E(x - \bar{x})^2 = 120$

\therefore Mean of Sampling distribution of means $= \mu_{\bar{x}} = \frac{\sum x}{n} = \frac{120}{10} = 12$

\therefore Variance of Sampling distribution of means is

$$\sigma_{\bar{x}}^2 = \frac{1}{n} \sum (x - \bar{x})^2 = \frac{1}{10} (120) = 12$$

$$\therefore \sigma_{\bar{x}} = \sqrt{12} = 3.464.$$

CHECK Mean of population $= \frac{\sum x}{N} = \frac{3+6+9+15+27}{5} = 12$

Variance of Population $= \frac{1}{N} \sum (x - \mu)^2$

$$= \frac{1}{5} [(3-12)^2 + (6-12)^2 + (9-12)^2 + (15-12)^2 + (27-12)^2]$$

$$= \frac{1}{5} [81 + 36 + 36 + 81 + 225] = \frac{370.2}{5} = 74.04$$

$$\therefore \sigma = \sqrt{74.04} = 8.602$$

Relation For finite population, $\mu_{\bar{x}} = \mu = 12$

$$\text{Variance } \sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \left(\frac{N-n}{N-1} \right) = \left(\frac{74.04}{4} \right) \frac{11.7}{3} = \frac{74.04}{12}$$

$$\sigma_{\bar{x}} = \sqrt{11.7} = 3.42$$

4. A normal population has a mean of 0.1 & S.D of 2.1. Find the probability that mean of the sample of size 900 will be negative.

Soln sample size = $n = 900$

$$\mu = \text{mean of population} = 0.1$$

$$\sigma = \text{S.D of Population} = 2.1$$

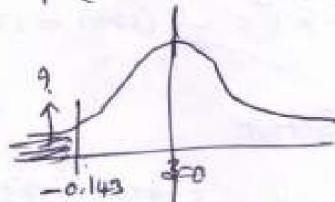
By using Central limit theorem, we have

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

we have to find $P(\bar{x} < 0)$.

$$\text{for } \bar{x} = 0 \Rightarrow z = \frac{-0.1}{2.1/\sqrt{900}} = -0.1 \cdot \frac{30}{2.1} = -0.143 - 1.43$$

$$\therefore P(\bar{x} < 0) = P(z < -0.143)$$



$$\therefore P(\bar{x} < 0) = P(z < -0.143)$$

$$= 0.5 - P(0 < z < 0.143)$$

$$= 0.5 - P(0 < z < 1.43)$$

$$= 0.5 - 0.4236$$

$$= 0.0764$$

5. A random sample of size 100 is taken from an infinite population having mean $\mu = 76$ & Variance $\sigma^2 = 256$. What is the probability that sample mean will be b/w 75 & 78.

Sol Given $n = 100$; $\mu = 76$; $\sigma^2 = 256$
 $\Rightarrow \sigma = \sqrt{256} = 16$

We have to find $P(75 < \bar{x} < 78)$

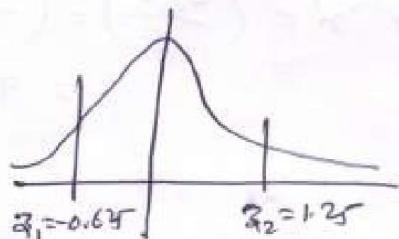
By using central limit theorem,

$$\text{for } \bar{x} = 75, z_1 = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{75 - 76}{16/\sqrt{10}} = -1 \left(\frac{10}{16}\right) = -0.625$$

$$\text{for } \bar{x} = 78, z_2 = \frac{78 - 76}{16/\sqrt{10}} = 2 \left(\frac{10}{16}\right) = 1.25$$

$$\therefore P(75 < \bar{x} < 78) = P(-0.625 < z < 1.25)$$

From normal table, we have



$$\begin{aligned} P(75 < \bar{x} < 78) &= P(-0.625 < z < 0) + P(0 < z < 1.25) \\ &= P(0 < z < 0.625) + P(0 < z < 1.25) \\ &= 0.2324 + 0.3944 \\ &\approx 0.6268 \\ &\approx 0.63 \end{aligned}$$

6. A sample of size 64 & mean 60 was taken from a population whose S.D is 10. Construct 95% confidence interval for the mean

Sol: Given $n=64$; $\bar{x}=60$; $\sigma=10$.

Confidence interval for the mean

$$= \left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$$

$$= \left(60 - \frac{10}{\sqrt{64}} (1.96), 60 + 1.966 \left(\frac{10}{8} \right) \right)$$

$$= (57.55, 62.45)$$

7. If maximum error is 5 & S.D of the population is 80 with 95% confidence, find the sample size

Sol: Given $E=5$; $\sigma=80$; $z_{\alpha/2}=1.96$

$$\therefore \text{Sample size } n = \left(\frac{z_{\alpha/2} \sigma}{E} \right)^2 = \left(\frac{1.96(80)}{5} \right)^2 \\ = 983.44$$

1. A random sample of size 64 is taken from a normal population with $\mu = 51.4$ and $\sigma = 6.8$. What is the probability that the mean of the sample will be
 a) exceed 52.9 b) fall between 50.5 & 52.3 c) be less than 50.6

SQ Given $n = 64$

Since Sample taken from Normal population

we use Central limit theorem.

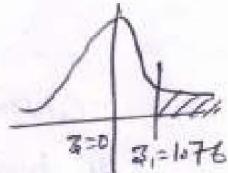
Also $\mu = 51.4$ and $\sigma = 6.8$

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

Now

a) we have to find $P(\bar{x} > 52.9)$

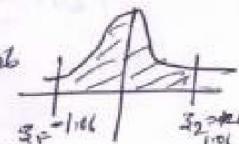
$$\text{For } \bar{x} = 52.9 \Rightarrow Z_1 = \frac{52.9 - 51.4}{6.8/\sqrt{64}} = 1.76$$



$$\begin{aligned} \therefore P(\bar{x} > 52.9) &= P(Z > 1.76) \\ &= 0.5 - P(0 < Z < 1.76) \quad (\text{from normal curve}) \\ &= 0.10392 \end{aligned}$$

b) we have to find $P(50.5 \leq \bar{x} \leq 52.3)$

$$\text{For } \bar{x}_1 = 50.5 \Rightarrow Z_1 = \frac{50.5 - 51.4}{6.8/\sqrt{64}} = -1.06$$



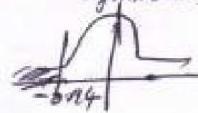
$$\bar{x}_2 = 52.3 \Rightarrow Z_2 = \frac{52.3 - 51.4}{6.8/\sqrt{64}} = 1.06$$

$$\therefore P(50.5 \leq \bar{x} \leq 52.3) = P(-1.06 < Z < 1.06)$$

$$\begin{aligned} &= 2P(0 < Z < 1.06) = 0.7108 \\ &\quad (\text{Since normal curve is symmetrical}) \end{aligned}$$

c) we have to find $P(\bar{x} < 50.6)$

$$\text{For } \bar{x} = 50.6 \Rightarrow Z_1 = \frac{50.6 - 51.4}{6.8/\sqrt{64}} = -0.94$$



$$\begin{aligned} P(\bar{x} < 50.6) &= P(Z < -0.94) \\ &= 0.5 - P(0 < Z < 0.94) = 0.5 - 0.3264 = 0.1736 \end{aligned}$$

2. A die is tossed 256 times and its turnup with an even digit 150 times. Is the die biased?

Sol? Probability of getting even digit in a single throw of a die (i.e 2 or 4 or 6) = $\frac{3}{6} = \frac{1}{2} = p$

$$\therefore q = 1 - p = 1 - \frac{1}{2} = \frac{1}{2}$$

$$\therefore \text{Mean} = np \approx 256\left(\frac{1}{2}\right) = 128 = \mu_{64}$$

$$\sigma^2 = \text{Variance} = npq = 256\left(\frac{1}{2}\right)^2 = \frac{256}{4} = 64$$

$$\Rightarrow \text{SD} = 8 = \sigma$$

$$\therefore \mu = 128, \sigma = 8$$

Given $x = \text{no of successes} = 150$

Null hypothesis H_0 : Die is unbiased

Alternative hypothesis H_1 : Die is biased

Level of significance: $\alpha = 0.05$

Test Statistic

$$Z_{\text{cal}} = \frac{x - \mu}{\sigma} = \frac{150 - 128}{8} = \frac{22}{8} = 2.75$$

Since at 5% level, $Z_{\text{tab}} = 1.96$

$\therefore Z_{\text{cal}} > Z_{\text{tab}} \Rightarrow \text{Null hypothesis rejected}$

\Rightarrow Alternative hypothesis accepted

\therefore The die is biased.

3. A sample of 400 items is taken from a population whose S.D is 10. The mean of the sample is 40. Test whether the sample has come from a population with mean 38. Also calculate 95% confidence interval for the population.

Solⁿ Given $n=400$; $\sigma=10$; $\bar{x}=40$

1. Null hypothesis $H_0: \mu=38$

2. Alternative hypothesis $H_1: \mu \neq 38$

3. Level of significance: $\alpha=0.05$

4. Test statistic

$$\bar{z}_{\text{cal}} = \frac{\bar{x}-\mu}{\sigma/\sqrt{n}} = \frac{40-38}{10/\sqrt{400}} = 4$$

since \bar{z}_{tab} at $\alpha=0.05 = 1.96$

$\bar{z}_{\text{cal}} > \bar{z}_{\text{tab}}$, we reject the null hypothesis
i.e. the sample is not taken from the population

with mean $\mu=38$.

95% Confidence Interval is given by

$$= \left(\bar{x} - 1.96 \left(\frac{\sigma}{\sqrt{n}} \right), \bar{x} + 1.96 \left(\frac{\sigma}{\sqrt{n}} \right) \right)$$

$$= \left(40 - 1.96 \left(\frac{10}{\sqrt{400}} \right), 40 + 1.96 \left(\frac{10}{\sqrt{400}} \right) \right)$$

$$= \left(40 - 0.98, 40 + 0.98 \right)$$

$$= (39.02, 40.98)$$

4. An ambulance service claims that it takes on the average less than 10 mins to reach its destination in emergency calls. A sample of 36 calls has a mean of 11 mins and the variance of 16 mins. Test the claim at 0.05 level of significance.

Soln Given $n=36$ (large sample)

$$\bar{x} = 11 \text{ mins}; S^2 = 16 \Rightarrow S = 4 \text{ mins}$$

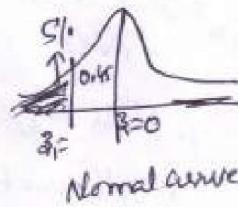
1. Null Hypothesis $H_0: \mu = 10 \text{ mins}$

2. Alternative Hypothesis $H_1: \mu < 10 \text{ mins}$

3. Level of significance: $\alpha = 0.05$

4. Test Statistic

$$Z = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{11 - 10}{4/\sqrt{36}} = \frac{6}{4} = 1.5$$



Normal curve

Since the problem is of one-tail test-

From normal curve

$$P(Z < Z_{tab}) = 0.45$$

$$\Rightarrow Z_{tab} = 1.645$$

$\therefore Z_{cal} < Z_{tab} \Rightarrow \text{Accept Null Hypothesis } H_0$

i. An ambulance service takes on average to reach its destination = 10 mins.

5. In 64 randomly selected hours of production, the mean & S.D of the number of acceptance pieces produced by an automatic stamping machine are 1.038 and 0.146. At the 0.05 level of significance does this enable us to reject the Null hypothesis $H_0: \mu = 1$ against the alternative hypothesis $H_1: \mu > 1$?

Sol: Given $n=64$, $\bar{x}=1.038$; $s=0.146$.

1. Null hypothesis $H_0: \mu=1$
2. Alternative hypothesis $H_1: \mu>1$ (one-tail test)
3. Level of significance: $\alpha=0.05$
4. Test statistic:

$$z = \frac{\bar{x}-\mu}{s/\sqrt{n}} = \frac{1.038-1}{0.146/\sqrt{64}} = 2.082$$

We know that $z_{\text{tab}}=1.645$ (for one-tailed test)

Since $z_{\text{cal}} > z_{\text{tab}}$, we reject H_0

$\therefore \mu > 1$
i.e. Mean of the population is greater than 1.

6. A company claims that its bulbs are superior to those of its main competitor. If a study showed that a sample of 40 of its bulbs have a mean life time of 647 hrs. of continuous use with a standard deviation of 27 hrs. Test the significance between the difference of two means at 5% level if a sample of 40 bulbs made by main competitor had a mean life

Ques: of 638 hrs of continuous use with SD of 31 hrs.

Sol: Given $n_1=40$; $\bar{x}_1=647$, $s_1=27$ hrs

Sol: Given $n_1=40$; $\bar{x}_1=647$, $s_1=27$ hrs
 $n_2=40$; $\bar{x}_2=638$, $s_2=31$ hrs

1. Null Hypothesis $H_0: \mu_1=\mu_2$

2. Alternative Hypothesis $H_1: \mu_1 > \mu_2$

3. Level of Significance: $\alpha = 0.05$

4. Test Statistic

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{647 - 628}{\sqrt{\frac{271^2}{40} + \frac{311^2}{40}}} = \frac{9}{\sqrt{\frac{729 + 961}{40}}}$$

$$\therefore Z_{\text{cal}} = \frac{9}{6.5} = 1.38$$

Since $Z_{\text{tab}} = 1.645$ (\therefore one tailed test)

$\because Z_{\text{cal}} < Z_{\text{tab}}$, we accept Null hypothesis H_0

i.e $\mu_1 = \mu_2$

Implies that the difference b/w two sample means is not significant.

7. In a certain factory there are two independent processes for manufacturing the same item. The average weight in a sample of 700 items produced from one process is found to be 250 gms with S.D of 30 gms while in a sample of 300 items from other process are 300 and 40. Is there significant difference b/w the means at 1% level.

Sol. Given $n_1 = 700$; $\bar{x}_1 = 250$ gms; $s_1 = 30$ gms

$n_2 = 300$; $\bar{x}_2 = 300$ gms; $s_2 = 40$ gms

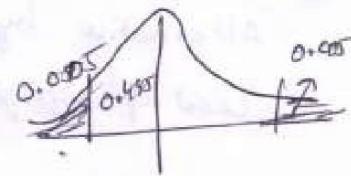
Let the Null hypothesis be $H_0: \mu_1 = \mu_2$

Alternative hypothesis $H_1: \mu_1 \neq \mu_2$ (Two tailed test)

Level of Significance: $\alpha = 0.01$

Test Statistic

$$\bar{z}_{\text{Cal}} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{250 - 300}{\sqrt{\frac{30^2}{700} + \frac{40^2}{300}}}$$



$$\Rightarrow \bar{z}_{\text{Cal}} = \frac{-50}{\sqrt{\frac{9}{7} + \frac{16}{3}}} = \frac{-50}{\sqrt{\frac{27+112}{21}}} = -50\sqrt{\frac{21}{139}} = -19.823$$

But \bar{z}_{tab} at 1% level of significance from

Normal Curve, $P(-\infty < \bar{z} < \bar{z}_{\text{tab}}) = 0.495$

$$\therefore \bar{z}_{\text{tab}} = 2.58$$

Since $|\bar{z}_{\text{Cal}}| > |\bar{z}_{\text{tab}}|$

we reject the Null hypothesis H_0

1. There is a significant difference between the means at 1% level of significance.

Ex. Samples of students were drawn from two universities and from their weights in kgs, Mean and S.D are calculated and shown below. Test the significance of the difference between the means.

	Mean	S.D	size
A	55	10	400
B	57	15	100

So. Given $n_1 = 400; \bar{x}_1 = 55; s_1 = 10$
 $n_2 = 100; \bar{x}_2 = 57; s_2 = 15$

Let Null hypothesis be there is no significant difference b/w the means i.e $H_0: \mu_1 = \mu_2$

Alternative hypothesis $H_1: \mu_1 \neq \mu_2$

Level of significance $\alpha = 0.01$ (10% level or
90% confidence)

Test statistic

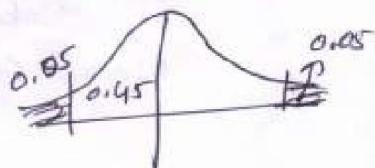
$$\bar{z}_{\text{cal}} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{55 - 57}{\sqrt{\frac{100}{400} + \frac{225}{100}}} = \frac{-2}{\sqrt{\frac{1}{4} + \frac{9}{4}}} = \frac{-2}{\sqrt{10}} = -1.96$$

0.45

Since 10% level of significance (two tailed test)

By using Normal curve, we get

$$P(0 < |z| < z_{\text{tab}}) = 0.45$$



$$\Rightarrow z_{\text{tab}} = 1.645$$

as $|z_{\text{cal}}| < z_{\text{tab}}$, we accept Null Hypothesis H_0

\Rightarrow There is no significant difference b/w the means
at 10% level of significance.

9. A manufacturer claimed that atleast 95% of the equipment which he supplied to a factory conformed to specifications. An examination of a sample of 200 pieces of equipment revealed that 18 were faulty. Test his claim at 5% level of significance. Also find its confidence interval at 5% level. [Confidence interval = $(p \pm 1.96 \sqrt{\frac{pq}{n}})$
 $= (0.8843, 0.9356)$]

Sol. Given the claim that atleast 95% of the equipment supplied to a factory is conformed to specifications. Use formula $(p - 2 \sqrt{\frac{pq}{n}}, p + 2 \sqrt{\frac{pq}{n}})$
 $=$ Confidence interval at 5% level

(16)

Also $n = 200$

$$\text{Proportion of faulty items in a sample} = \frac{18}{200}$$

$$\Rightarrow \text{Proportion up to specifications in a sample} = 1 - \frac{18}{200} \\ = \frac{182}{200}$$

$$\therefore p = \frac{91}{100} = 0.91$$

Let Null hypothesis be $H_0: P = 0.95$ Alternative hypothesis $H_1: P \neq 0.95$ (One-tailed test)Level of significance: $\alpha = 0.05$

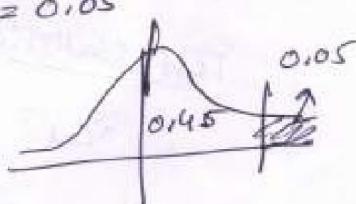
Test Statistic

$$Z_{\text{Cal}} = \frac{p - P}{\sqrt{\frac{P(1-P)}{n}}}$$

$$\text{Here } \alpha = 1 - P = 1 - 0.95 = 0.05$$

$$\therefore Z_{\text{Cal}} = \frac{0.91 - 0.95}{\sqrt{\frac{(0.95)(0.05)}{200}}}$$

$$\Rightarrow Z_{\text{Cal}} = \frac{-0.04}{\sqrt{\frac{0.049}{200}}} = \frac{-0.04}{\sqrt{0.000245}} = \frac{-0.04}{0.0156} = -2.564$$



At 5%, level of significance (one-tailed)

By Normal curve we have

$$Z_{\text{tab}} = 1.645$$

 $\because |Z_{\text{Cal}}| > Z_{\text{tab}}$ we reject H_0

$$\Rightarrow P \neq 0.95$$

 $\Rightarrow P < 0.95 \Rightarrow \text{Manufacture's claim is rejected.}$

10. In a big city 325 men out of 600 men were found to be smokers. Does this information support the conclusion that the majority of men in this city are smokers?

Sol. Given $n = 600$

$$\text{Proportion of smokers in sample} = \frac{325}{600} = \frac{13}{24} = 0.5417$$

Let the Null hypothesis be $H_0: P \geq \frac{1}{2}$

i.e. majority of men in a city are ^{not} smokers.

Alternative Hypothesis $H_1: P > \frac{1}{2}$ (one tailed)

i.e. Majority of men in a city are smokers.

Level of significance: $\alpha = 0.05$

$$\text{Test statistic} = \frac{\hat{P} - P}{\sqrt{\frac{PQ}{n}}} \quad \text{where } \alpha = \frac{1}{2}$$

$$= \frac{0.5417 - 0.5}{\sqrt{\frac{(0.5)(0.5)}{600}}} = \frac{0.0417}{\sqrt{\frac{0.5}{600}}} = 2.04$$

At 5% level of significance,

$$\hat{Z}_{\text{tab}} = 1.645$$

$\hat{Z}_{\text{cal}} > \hat{Z}_{\text{tab}} \Rightarrow$ we reject H_0 and we conclude

that the Majority of men in the city are smokers.

11. Experience had shown that 20% of a manufactured product is of the top quality. In one day's production of 400 articles only 50 are of top quality. Test the hypothesis at 0.05 level.

Sol. Given Proportion of top quality in population = $P = \frac{20}{100} = 0.2$.

$n = 400$; Proportion of top quality in sample is

$$p = \frac{50}{400} = 0.125$$

1. let Null hypothesis be $H_0: P = 0.2$. Then $\alpha = 0.8$

2. Alternative hypothesis $H_1: P \neq 0.2$

3. Level of Significance: $\alpha = 0.05$

4. Test statistic

$$\hat{\sigma}_{\text{cal}} = \sqrt{\frac{p - P}{n}} = \sqrt{\frac{0.125 - 0.2}{400}} = \sqrt{\frac{-0.075}{0.02}} = -3.75$$

Since $\hat{\sigma}_{\text{tab}} = 1.96$ at $\alpha = 0.05$, we have

$\therefore |\hat{\sigma}_{\text{cal}}| > \hat{\sigma}_{\text{tab}}$ implies that Reject the

Null hypothesis, \therefore Reject

So, we conclude that experience claim is not correct.

12. In two large populations, there are 30% and 25% respectively of fair haired colors people. Is this difference likely to be hidden in samples of 1200 and 900 respectively from the two populations.

Sol: Given proportion of fair haired in first

$$\text{population} = \frac{30}{100} = 0.3 = P_1$$

Proportion of fair haired in second population

$$= \frac{25}{100} = \frac{1}{4} = 0.25 = P_2$$

Given $n_1 = 1200$ & $n_2 = 900$

let Null Hypothesis H_0 : Assume that the sample proportions are equal i.e. the difference in population proportions is likely to be hidden in sampling i.e. $P_1 = P_2$

Alternative Hypothesis H_1 : $P_1 \neq P_2$

Level of significance: $\alpha = 0.05$

Test statistic

$$Z_{\text{cal}} = \frac{P_1 - P_2}{\sqrt{\frac{P_1(1-P_1)}{n_1} + \frac{P_2(1-P_2)}{n_2}}} = \frac{0.3 - 0.25}{\sqrt{\frac{0.3(0.7)}{1200} + \frac{0.25(0.75)}{900}}} = 2.56$$

We know Z_{tab} at 5%, level = 1.96

since $Z_{\text{cal}} > Z_{\text{tab}}$, we reject H_0 and we conclude that the difference in population proportions is unlikely that the real difference will be hidden.

Q. 100 articles from a factory are examined and 10 are found to be defective. 500 similar articles from a second factory are found to be 15 defective. Test the significance b/w the difference of the two proportions at 1% level.

Sol. Given $n_1 = 100$; No of articles defective in first sample = 10
 \therefore proportion of defective items in first sample $= P_1 = 0.1$
 $n_2 = 500$; proportion of defective items in second sample $= \frac{15}{500} = 0.03 = P_2$

Null Hypothesis $H_0: P_1 = P_2$

Alternative Hypothesis $H_1: P_1 \neq P_2$

Level of significance: $\alpha = 0.01$

Test statistic

$$\bar{Z}_{\text{cal}} = \frac{P_1 - P_2}{\sqrt{\frac{P(1-P)}{n_1 + n_2}}} \quad (\text{Method of pooling})$$

$$\text{Where } p = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2} = \frac{100(0.1) + 500(0.03)}{600} = \frac{25}{600} = 0.0416$$

$$q = 1 - p = 1 - 0.0416 = 0.958$$

$$\therefore \bar{Z}_{\text{cal}} = \frac{0.1 - 0.03}{\sqrt{(0.0416)(0.958)\left(\frac{1}{100} + \frac{1}{500}\right)}} = \frac{0.07}{0.022} = 3.18$$

At 1% level, $Z_{\text{tab}} = 2.58$

$\therefore \bar{Z}_{\text{cal}} > Z_{\text{tab}}$ we reject $H_0 \Rightarrow$ There is significant difference b/w two proportions

14. In an investigation on the machine performance the following results are obtained.

	Inspected	Defects
Mach A	375	17
B	450	22

Test whether there is any significant performance of two machines at $\alpha = 0.05$.

SOL Given $n_1 = 375$; $n_2 = 450$; $p_1 = \frac{17}{375}$; $p_2 = \frac{22}{450} = 0.049$
 $= 0.0453$

Null Hypothesis $H_0: P_1 = P_2$

Alternative Hypothesis $H_1: P_1 \neq P_2$

Level of significance: $\alpha = 0.05$

Test Statistic

$$Z_{\text{cal}} = \frac{P_1 - P_2}{\sqrt{pq(\frac{1}{n_1} + \frac{1}{n_2})}} \quad (\text{Method of Pooling})$$

$\uparrow \text{if } p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$

$$= \frac{0.0453 - 0.049}{\sqrt{(0.0472)(0.953)(\frac{1}{375} + \frac{1}{450})}} = \frac{-0.257}{\sqrt{0.015}} = -0.257$$

$$\therefore |Z_{\text{cal}}| = 0.257$$

Since $|Z_{\text{cal}}| < Z_{\text{tab}} = 1.96$, we accept H_0

\Rightarrow There is no significant performance of two machines at 5% level of significance

Exact Sampling Distributions Unit - 4 . (Small Samples) ①

When the sample is small, the Sampling Distribution in many cases may not be normal. Here the test statistics will change.

Degrees of Freedom : denoted by ν . It is the no. of values in a set of data which may be assigned arbitrarily or, it refers to the no. of independent constraints in a set of data.

In general, the no. of degrees of freedom = to the total no. of observations ~~-~~ - no. of independent constraints imposed on the observations.

* In brief, The no. of independent variates which make up the statistic is known as degrees of freedom (d.f.)

t-distribution (OR) Student's t-Dist.

$$t = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \quad \text{with } \nu = n - 1, \text{ d.f.}$$

where \bar{x} - sample mean

μ - Population "

n - sample size

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (x_i - \bar{x})^2$$

Prob. density fn: $f(t) = y_0 \left(1 + \frac{t^2}{2}\right)^{-\frac{n+1}{2}}$ where

y_0 is a constant got by $\int_{-\infty}^{\infty} f(t) dt = 1$.

t-distr. is extensively used in hypothesis about one mean, or about equality of two means when σ is unknown.

Applications of t-Distribution

- (1) To test the significance of the sample mean, when population variance (σ^2) is not given.
- (2) To test the significance of the mean of the sample, i.e. to test if the sample mean differs significantly from the population mean (μ).
- (3) To test the significance of the difference between two sample means or to compare 2 samples.
- (4) To test the significance of an observed sample correlation coefficient & sample regression coefficient.

Chi-Squared (χ^2) Distribution:

It is a conti prob. distri of a conti r.v. X with prob. density fn given by

$$f(x) = \begin{cases} \frac{1}{2^{k/2} \Gamma(k/2)} x^{(k/2)-1} e^{-x/2}, & \text{for } x > 0 \\ 0, & \text{otherwise} \end{cases}$$

χ^2 -distri. was extensively used as a measure of goodness of fit & to test the independence of attributes.

Note (1) χ^2 values from 0 to ∞ .

(2) If χ_1^2 & χ_2^2 are 2-independent distributions with k_1 & k_2 dof, then $\chi_1^2 + \chi_2^2$ will be chi-sqrd distri. with $(k_1 + k_2)$ dof. i.e. It is additive.

* It is used in sampling distribution analysis of variance, mainly it is used as a measure of goodness of fit & in analysis of $r \times c$ tables.

$$\chi^2 = \frac{(n-1)s^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} \quad \{ s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \}$$

χ^2 is a value of a r.v having χ^2 distribution with $df = n-1$.

Note: (i) Exactly 95% of χ^2 distri: lies bet' $\chi^2_{0.975}$ & $\chi^2_{0.025}$

(ii) & when σ^2 is too small

~~(iii)~~ χ^2 falls to rt. of $\chi^2_{0.025}$ & χ^2 falls to the left of

(iii) when σ^2 is too large

$\chi^2_{0.975}$. Thus when σ^2 is correct, χ^2 values are

to the left of $\chi^2_{0.975}$ or to the right of $\chi^2_{0.025}$.

F-distribution (S.D. of the ratio of 2 sample
It is another imp conti. prob. variances).
distri.

$$F = \frac{\frac{s_1^2}{\sigma_1^2}}{\frac{s_2^2}{\sigma_2^2}} = \frac{s_1^2}{s_2^2}$$

where s_1^2, s_2^2 are variances of 2 random samples of sizes n_1 & n_2 with variances σ_1^2 & σ_2^2 . which follows f-distri. with $df_1 = n_1 - 1$ & $df_2 = n_2 - 1$.

Note (1) F determines whether the ratio of 2 sample variances s_1^2, s_2^2 is too small or too large.

(2) When F is close to 1, s_1, s_2 are almost same.

(3) F is always a +ve no.

San-Diss if F is of the form $f(F) = K \frac{F^{(\beta_1-2)/2}}{(\nabla_1 F + \nabla_2)^{(\beta_1+\beta_2)/2}}$

where K is determined by $\int_0^{\infty} f(F)dF = 1.$

t -distribution

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad \text{where } s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

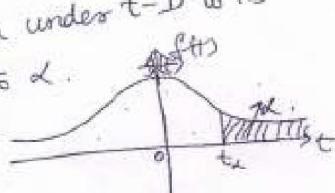
$$\text{p.d.f} \quad f(t) = y_0 \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

where $\nu = n-1$ dof
 y_0 is a constant got by $\int_{-\infty}^{\infty} f(t) dt = 1$.

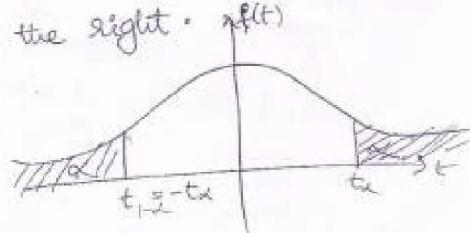
Note: The mean of std ND as well as t-D is zero, but variance of t-D depends on parameter ν .

- t-D with ν dof approaches std ND as $\nu \rightarrow \infty$

- t_α denotes area under t-D to its right is equal to α .



* $t_{1-\alpha} = -t_\alpha$
 i.e. the t-value ~~leaving~~ leaving an area of $1-\alpha$ to the right & \therefore an area α to its left = -ve t-value which leads an area α in the right.



No
 t-D is extensively used in hypothesis about one mean, or about equality of 2 means when σ is unknown.

Applications:

- 1) To test significance of sample mean, when pop. var. is not given.
- 2) To test significance of the mean of the sample i.e. to test if sample mean differs significantly from pop. mean.
- 3) To Test significance of difference between 2 sample means or to compare 2 samples.
- 4) -

Chi-Squared (χ^2) Distribution:

p.d.f is given by

$$f(x) = \begin{cases} \frac{1}{2^{n/2} \Gamma(n/2)} x^{(n/2)-1} e^{-x/2}, & \text{for } x > 0 \\ 0, & \text{otherwise} \end{cases}$$

Properties:

(1) χ^2 is not symmetrical, lies entirely in I Quadrant & hence not a normal curve, $\because \chi^2$ varies from 0 to ∞ .

(2) It depends only on n .

— Mean = n , Variance = $2n$

$$\chi^2 = \frac{(n-1)s^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} \quad \left. \begin{array}{l} s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \end{array} \right\}$$

Note: Exactly 95% of χ^2 -D lies between

$$\chi^2_{0.975} \text{ & } \chi^2_{0.025}$$

- (i) When σ^2 is too small, χ^2 falls to the right of $\chi^2_{0.025}$
- (ii) When σ^2 is too large, χ^2 falls to the left of $\chi^2_{0.975}$.

When σ^2 is correct χ^2 -values are to the left of $\chi^2_{0.975}$ or to the right of $\chi^2_{0.025}$.

Applications:

- (1) To test goodness of fit.
- (2) " " independence of attributes.
- (3) " " homogeneity of independent estimation of population variance (i.e.)
- (4) To test homogeneity of independent estimation of population correlation coefficient.

F-Distribution

(Sampling Distr. of ratio of 2 Sample Variances)
 To determine whether the 2 samples come from 2 populations having equal variances.

$$F = \frac{s_1^2 / \sigma_1^2}{s_2^2 / \sigma_2^2} = \frac{\sigma_2^2 s_1^2}{\sigma_1^2 s_2^2}$$

with $v_1 = n_1 - 1$ & $v_2 = n_2 - 1$ d.f

Note: F-D can be used to test equality of several population means, comparing sample variances & analysis of variance.

$F = \frac{s_1^2}{s_2^2}$ under the assumption that 2 normal populations have the same variance. i.e. $\sigma_1^2 = \sigma_2^2$.

- ★ F determines whether the ratio of 2 sample variances s_1^2 & s_2^2 is too small or too large.
 When F is close to 1, the 2 sample variances s_1^2 & s_2^2 are almost same.

note: In practice, it is customary, to take the larger sample variance as the numerator.
 - F is always a +ve no.

The sampling Distribution of F is of the form $f(F) = K \frac{F^{(v_2-2)/2}}{(v_1 F + v_2)^{(v_1+v_2)/2}}$ where
 K is determined by $\int_0^\infty f(F) dF = 1$.

Properties .

- 1) F-D curve lies entirely in first-quadrant.
- 2) The F-curve depends not only on v_1, v_2 but also on the order in which they are used .
- 3) $F_{1-\alpha}(v_1, v_2) = \frac{1}{F_x(v_2; v_1)}$ where
 $F_x(v_1, v_2)$ is the value of F with v_1, v_2 dof \Rightarrow area under F-D curve to the right of F_x is α .
- 4) mode of F-D is less than unity .

Sampling Distribution of Proportions: (SDP)

Let 'p' be prob. of occurrence of an event (called success) &

$q = 1-p$ is prob. of non-occurrence (called failure)

Draw all possible samples of size 'n' from an infinite population .

Compute P proportion of success for each of these samples .

Mean μ_p & Variance σ_p^2 of ~~SDP~~ are

$$\mu_p = p , \sigma_p^2 = \frac{pq}{n} = \frac{p(1-p)}{n}$$

Note: For finite population (with replacement) of size N

$$\mu_p = p \Rightarrow \sigma_p^2 = \frac{pq}{n} \left(\frac{N-n}{N-1} \right)$$

Not
Needed

Unit IVTest of Significance for Small Samples

On the basis of sample results, the tests of significance enable us to decide, if

- (i) the deviation b/w Observed Sample statistic & hypothetical parameter value is significant.
- (ii) the deviation b/w 2 sample statistics is significant.

Some imp. tests for small samples are :

- (a) Student's 't'-test
- (b) F-test
- (c) χ^2 -test

Student's 't' test :

Assumptions of 't'-test :

- (i) Sample size , $n < 30$
- (ii) The parent population from which sample is drawn is Normal .
- (iii) The population standard deviation (may or maynot be known) is unknown .
- (iv) The sample observations are independent . i.e sample is random .

Uses : This test is used

- (1) to test for a specified mean
- (2) to test for equality of 2 means (of 2 independent samples drawn from 2 normal populations), S.D of the populations being unknown
- (3) to test the significance of difference b/w the means of paired data .

let \bar{x} = mean of a sample
 n = size of the "
 s = S.D of " "
 μ = Mean of the population supposed to be normal.
 student's t is defined by the statistic

$$t = \frac{\bar{x} - \mu}{S/\sqrt{n-1}}$$

case(i) when SD of sample i.e 's' is given directly then

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n-1}}$$

case(ii) when 's' is not given but samples are given

then calculate 's' from $s^2 = \frac{\sum (x_i - \bar{x})^2}{n}$

& then use in the formula above.

Note: Also if $S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$ then $t = \frac{\bar{x} - \mu}{S/\sqrt{n}}$

where S^2 is called unbiased estimate of population variance σ^2 .

In general while solving problems we ^{use} ~~use~~ σ for s

* The Confidence or Fiducial limits for μ are

$$\left[\bar{x} \pm t_{0.05} \cdot \frac{s}{\sqrt{n}} \right] \text{ at } 5\% \text{ Level of significance}$$

& for $(n-1)$ degrees of freedom.

$$\text{In general, } \bar{x} \pm t_{\alpha} \frac{s}{\sqrt{n}}$$

: Students 't' Test for Single Mean .

To test the hypothesis that population mean μ has a specified value μ_0 when population S.D σ is not known.

$$H_0: \mu = \mu_0$$

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

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n-1}} \quad \text{where 's' is sample S.D with}$$

$\Rightarrow (n-1)$ degrees of freedom.

$\alpha = 0.05$ (ie 95% confidence limits).

Find table value of 't' at α level of significance.

If calculated value of $|t| >$ table value of t

\Rightarrow Reject H_0

If $|t| < t_{\text{table value}}$ \Rightarrow Accept H_0 .

Note : For a 2-tailed test , value of $\frac{\alpha}{2}$ is taken for α .

Do problems .

Student's t-test for Difference of Means

Let μ_1, μ_2 be means of two Normal Populations

let n_1, n_2 be the sample sizes from the two populations.

Let \bar{x}, \bar{y} be the means of these 2 independent samples whose S.D.s are s_1, s_2 respectively.

To Test whether the 2 population means are equal.
(i.e. whether the difference $\mu_1 - \mu_2$ is significant).

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

$$t = \frac{\bar{x} - \bar{y}}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\Rightarrow (n_1 + n_2 - 2) \text{ degrees of freedom.}$$

$$\text{Here } S = \sqrt{\frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2 - 2}}$$

Note: If s_1, s_2 are not given then

$$S^2 = \frac{1}{n_1 + n_2 - 2} \left[\sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2 \right] \text{ where}$$

$$\bar{x} = \frac{1}{n_1} \sum_{i=1}^{n_1} x_i; \quad \bar{y} = \frac{1}{n_2} \sum_{i=1}^{n_2} y_i$$

Obtain table value of 't' for '2' d.f.

if $|t| < t_\alpha \Rightarrow$ Accept H_0
otherwise Reject H_0 .

* Confidence limits for difference of 2 population mean
are $(\bar{x} - \bar{y}) \pm t_\alpha S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$

Note: Same procedure is followed whether samples are taken from different populations or same.

Paired Sample t-Test:

Suppose a business concern is interested to know whether a particular media of promoting sales of a product is really effective or not. Here we have to test whether the average sales "before" and "after" the sales promotion are equal.

Thus the particular media is the experimental unit & the 2 populations are average sales "before" and "after". Hence the two samples are not independent.

These are paired observations which arise in many practical situations where each homogeneous experimental unit receives both population conditions.

If $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be pairs of sales data "before" & "after" the sales promotion in a business concern, paired t-Test is applied to test the significance of the difference of the two situations.

Let $d_i = x_i - y_i$ (or $y_i - x_i$) for $i = 1, 2, 3, \dots, n$

$$H_0: \mu_1 = \mu_2 \text{ (i.e } \mu = 0\text{)}$$

$$H_1: \mu_1 \neq \mu_2$$

The test statistic for 'n' paired observations (which are dependent) by taking the differences d_1, d_2, \dots, d_n of the paired data.

$$t = \frac{\bar{d}}{s/\sqrt{n}} \quad \left. \begin{array}{l} \text{if } \mu = 0 \\ \text{otherwise } t = \frac{\bar{d} - \mu}{s/\sqrt{n}} \end{array} \right\}$$

$$\bar{d} = \frac{1}{n} \sum d_i, \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (d_i - \bar{d})^2$$

α = Level of Significance (1.0, 5).

$\nu = n-1$ d.f.

if H_1 \leftarrow Accept H_0 otherwise Reject H_0

Snedecor's F-test.

If variances of the populations are not equal, a significant difference in the means may arise. Hence before we apply t-test for significance of difference of two means, we have to test for equality of population variances using F-test.

$$\text{the test statistic } F = \frac{\text{Greater Variance}}{\text{Smaller Variance}}$$

when F is close to 1, the 2 sample variances (say s_1^2, s_2^2) are nearly same.

$F_\alpha(\nu_1, \nu_2)$ is value of F with ν_1, ν_2 d.f. such that the area under F-distribution to the right of F_α is α . Clearly value of F at 5% significance is lower than at 1%.

To test the hypothesis that the two population variances σ_1^2 & σ_2^2 are equal:

Let two independent random samples of sizes n_1, n_2 be drawn from 2 normal populations.

The estimates of σ_1^2, σ_2^2 are given by

$$S_1^2 = \frac{n_1 s_1^2}{n_1 - 1} = \frac{\sum (x_i - \bar{x})^2}{n_1 - 1}$$

$$S_2^2 = \frac{n_2 s_2^2}{n_2 - 1} = \frac{\sum (y_i - \bar{y})^2}{n_2 - 1}$$

where s_1^2, s_2^2 are variances of the two samples.

$$F = \frac{S_1^2}{S_2^2} \quad (\text{or } \frac{S_2^2}{S_1^2}) \text{ according as } S_1^2 > S_2^2 \text{ (or } S_2^2 > S_1^2)$$

$\nu_1 = n_1 - 1, \nu_2 = n_2 - 1$ are degrees of freedom

if $F_{\text{calculated value}} < F_{\text{tabulated value}}$ we accept H_0 &

conclude σ_1^2, σ_2^2 are equal . Otherwise Reject H_0 .

(Do problems)

Chi-Square Test : χ^2 -Test.

χ^2 Test is used to test whether differences b/w Observed & expected frequencies are significant. It is mainly used

- (1) To test the goodness of fit
- (2) To test the independence of attributes.
- (3) To test if the population has a specified value of the variance σ^2 .

In general, the results obtained in an experiment do not agree exactly with the theoretical results. The magnitude of discrepancy b/w the theory & observation is given by the quantity χ^2 (a Greek letter, read as Chi-square).

If $\chi^2 = 0$, observed & expected frequencies coincide completely.

As χ^2 value increases, the discrepancy b/w the two increases.

Thus χ^2 affords a measure of correspondence b/w theory & observation.

Note: If the data is given in a series of 'n' no.s then
 $d.f = n-1$.

In case of Binomial Distribution : $d.f = n-1$

" " Poisson " : $d.f = n-2$

" " Normal " : $d.f = n-3$.

If O_i ($i=1, 2, \dots, n$) is set of Observed (experimental)

frequencies, &
 E_i ($i=1, 2, \dots, n$) is corresponding set of expected (theoretical)

frequencies, then

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

$$\boxed{\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}}$$

with $d = n-1$ d.f.

χ^2 Test as Test of Goodness of Fit

Conditions of Validity:

The following conditions should be satisfied before χ^2 Test is applied.

- (i) Sample Observations should be independent.
- (ii) Total frequency (ie N) is large. ie $N > 50$.
- (iii) The constraints q on the cell frequencies, if any, are linear.
- (iv) No theoretical (or expected) frequency should be less than 10.

If they occur, then the difficulty is overcome by grouping 2 or more classes together before finding $(O-E)$.
Also the d.f is determined with no. of classes after regrouping.

H_0 : There is no significant difference b/w Observed values & expected values.

H_1 : The difference is significant.

$$\chi^2 = \sum_{i=1}^{n} \left[\frac{(O_i - E_i)^2}{E_i} \right]$$

$$d.f = n-1$$

$$d.o.f = L.O.S.$$

If Calculated $\chi^2 <$ Tabulated $\chi^2 \Rightarrow$ Accept H_0
Otherwise Reject H_0 .

Do problems

Chi-Square Test for Independence Of Attributes :

'Attribute' literally means a quality or characteristic.
 e.g.: drinking, smoking, blindness, honesty, beauty etc.
 An Attribute may be marked by its presence (position) or absence in a no. of a given population.
 Let Observed frequencies be classified according to two attributes say A, B & the frequencies O_i in the different categories be shown in a two-way table called contingency table.

On the basis of cell frequencies we have to Test whether the 2 attributes are independent or not.

The Expected frequencies (E_i) of any cell is

$$E_i = \frac{\text{Row total} \times \text{Column total}}{\text{Grand total}}$$

i.e. If

A	a	b	atb	ab
B	c	d	ctd	cb
	atc	btd	$N=atb+ctd$	

 Expected frequencies are:

$$E(a) = \frac{(a+c)(a+b)}{N}$$

$$E(b) = \frac{(b+d)(a+b)}{N}; E(c) = \frac{(a+c)(c+d)}{N}; E(d) = \frac{(b+d)(c+d)}{N}$$

The table is

O_i	E_i	$O_i - E_i$	$\frac{(O_i - E_i)^2}{E_i}$	$\sum_{i=1}^r \frac{(O_i - E_i)^2}{E_i}$
a	$E(a)$	$a - E(a)$	-	-
b	$E(b)$	$b - E(b)$	-	-
c	$E(c)$	$c - E(c)$	-	-
d	$E(d)$	$d - E(d)$	-	-

Degrees of freedom: $\gamma = (\text{No. of rows} - 1) \times (\text{No. of columns} - 1)$
 $\Rightarrow \gamma = (r-1)(c-1)$

Unit-IV
 Problems on Exact
Sampling Distribution
(Small Samples)

① 13

1. A sample of 26 bulbs gives a mean life of 990 hours with a S.D of 20 hrs. The manufacturer claims that mean life of bulbs is 1000 hrs. Is the sample not up to the standard.

Sol. Given $n=26$; $\bar{x}=990$, $s=20$

let Null hypothesis $H_0: \mu=1000$ hrs.

Alternative Hypothesis $H_1: \mu < 1000$

Level of Significance: $\alpha=0.05 + d.o.f(n-1)$

Test Statistic

$$t_{\text{cal}} = \frac{\bar{x}-\mu}{s/\sqrt{n-1}} = \frac{990-1000}{20/\sqrt{25}} = \frac{-10}{4} = -2.5$$

But t_{tab} at $\alpha=0.05$ with 25 d.o.f = 1.708

$\therefore |t_{\text{cal}}| > t_{\text{tab}}$; we reject the Null Hypothesis H_0

i.e $\mu \neq 1000 \Rightarrow$ The manufacturer claim Mean of bulbs 1000 hrs is not correct.

i.e the sample is not up to the standard.

2. A sample of 25 members has a mean 67 and S.D 5.2
 Is this sample has been taken from a large population
 of mean 70?

Sol. Given $n=25$, $\bar{x}=67$; $s=5.2$

let Null Hypothesis $H_0: \mu=70$

Alternative Hypothesis $H_1: \mu \neq 70$

level of significance: $\alpha = 0.05$ with 94 d.o.f

Test statistic

$$t_{\text{cal}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n-1}} = \frac{67 - 70}{5.4/\sqrt{95-1}} = \frac{-3}{1.6614} = -1.82.$$

$$\therefore |t_{\text{cal}}| = 1.82 > 2.06$$

we reject the Null hypothesis and concluding that the sample has not taken from population of mean μ_0 .

3. A random sample of 10 boys had the following I.Q's: 70, 120, 110, 101, 88, 83, 95, 98, 107 and 100.

a) Do these data support the assumption of a population mean I.Q of 100.

b) find a reasonable range in which most of the means

I.Q values of samples of 10 boys?

Given $n = 10$
Here S.D and Mean of the Sample are not given

we have to find Mean & S.D of sample

$$\text{So, Mean} = \frac{70 + 120 + 110 + 101 + 88 + 83 + 95 + 98 + 107 + 100}{10} = \frac{972}{10} = 97.2$$

$$\therefore \bar{x} = 97.2$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^{10} (x_i - \bar{x})^2 = \frac{1}{9} [(70-97.2)^2 + (120-97.2)^2 + \dots + (100-97.2)^2] \\ = \frac{1833.6}{9} = 203.73$$

$$\Rightarrow S = \sqrt{203.73} = 14.27.$$

Let Null Hypothesis $H_0: \mu = 100$

Alternative Hypothesis $H_1: \mu \neq 100$ (Two tailed test)

Level of significance: $\alpha = 0.05$ with 9 d.o.f

Test Statistic

Q 14

$$t_{\text{cal}} = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{97.2 - 100}{14.27/\sqrt{10}} = -0.62$$

$$\therefore |t_{\text{cal}}| = 0.62$$

Since $|t_{\text{cal}}|$ with 9 d.o.f and $\alpha = 0.05 \rightarrow 2.26$

we accept Null Hypothesis.

\therefore the assumption of mean 100 of 100 in the population is correct (supportive).

b) The 95% Confidence limits = $(\bar{x} - t_{0.05 \text{ with 9 d.o.f}} \cdot (S/\sqrt{n}),$

$$\bar{x} + t_{0.05} (S/\sqrt{n})] = (97 \pm 2.26(14.27/\sqrt{10})) = 87 \text{ and } 107.4.$$

4. The heights of 10 males of a given locality are found to be 70, 67, 62, 68, 61, 68, 70, 64, 64, 66 inches. Is it reasonable to believe that the average height is greater than 64 inches? Test at 1% level of significance

Sol. Given $n=10$; Mean = $\frac{\sum x}{n} = \bar{x}$; $S^2 = \frac{1}{n-1} \sum (x-\bar{x})^2$

x	$x - \bar{x}$	$(x - \bar{x})^2$
70	4	16
67	1	1
62	-4	16
68	2	4
61	-5	25
68	2	4
70	4	16
64	-2	4
64	-2	4
66	0	0
660		90
Σx		$\Sigma (x - \bar{x})^2$

$$\therefore \bar{x} = 66$$

$$\text{and } S^2 = \frac{1}{9} [90] = 10$$

$$S = \sqrt{10} = 3.16$$

1. Null Hypothesis $H_0: \mu = 64$
2. Alternative Hypothesis $H_1: \mu > 64$
3. Level of Significance: $\alpha = 0.01$ with 9 d.o.f
4. Test statistic

$$t_{\text{cal}} = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{66 - 64}{3.16/\sqrt{10}} = \frac{2\sqrt{10}}{3.16} = 2.0714$$

Here t_{tab} at 0.01 level of significance with 9 d.o.f
 $= 2.821$

Since $t_{\text{cal}} < t_{\text{tab}}$ we accept null hypothesis
 \therefore At 1% level of significance $\mu = 64$.

5. Two horses A and B were tested according to the time (in sec) to run a particular track with the following results.

Horse A	28	30	32	33	33	29	34
Horse B	29	30	30	24	27	29	-

Test whether the two horses have the same running capacity.

Sol: Given $n_1 = 7$; $n_2 = 6$
Degree of freedom $= n_1 + n_2 - 2 = 7+6-2 = 11$ d.o.f
Mean of series A = $\bar{x} = \frac{\sum x}{n_1} = \frac{28+30+32+66+29+34}{7} = 31.286$
Mean of series B = $\bar{y} = \frac{\sum y}{n_2} = \frac{29+60+24+27+29}{6} = 28.16$
Now $S^2 = \frac{1}{n_1+n_2-2} [\sum (x-\bar{x})^2 + \sum (y-\bar{y})^2]$

x	$\frac{(x-\bar{x})}{\bar{x}-31.286}$	$(x-\bar{x})^2$	y	$(y-\bar{y})$	$(y-\bar{y})^2$
28	-3.286	10.8	29	0.84	0.7056
30	-1.286	1.6528	30	1.84	3.2856
32	0.714	0.51	30	1.84	3.2856
33	1.714	2.94	24	-4.16	17.3056
33	1.714	2.94	27	-1.16	1.3456
33	1.714	5.226	29	0.84	0.7056
29	-2.286	7.366	-	-	
34	2.714	31.4358	169		26.8336
219					

$$\therefore \sum (x-\bar{x})^2 = 31.4358 \text{ and } \sum (y-\bar{y})^2 = 26.8336$$

$$\therefore S^2 = \frac{1}{11} [31.4358 + 26.8336] = \frac{58.2694}{11} = 5.23$$

$$\therefore S = \sqrt{5.23} = 2.3$$

1. Null Hypothesis $H_0: \mu_1 = \mu_2$

2. Alternative Hypothesis $H_1: \mu_1 \neq \mu_2$

3. Level of Significance: $\alpha = 0.05$ with 1 d.o.f

4. Test Statistic

$$t_{cal} = \frac{\bar{x} - \bar{y}}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{31.286 - 28.16}{2.3 \sqrt{\frac{1}{7} + \frac{1}{6}}} = 2.443$$

Here $t_{tab} = 2.201$

since $t_{cal} > t_{tab}$, we reject the Null Hyp H_0

\Rightarrow There is the significant difference in the running capacity of horses A and B.

6. find the Maximum difference that we can expect with probability 0.95 b/w the means of samples of sizes 10 and 12 from a normal population if their S.Ds are found to be 2 and 3 respectively.

Given $n_1 = 10$ and $n_2 = 12$

$$S_1 = 2, \quad S_2 = 3$$

$$\text{then } S^2 = \frac{n_1 S_1^2 + n_2 S_2^2}{n_1 + n_2 - 2} = \frac{10(4) + 12(9)}{20} = \frac{40 + 108}{20} = \frac{148}{20} = 7.4$$

$$\therefore S = \sqrt{7.4} = 2.72.$$

Test statistic

$$t = \frac{\bar{x} - \bar{y}}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Here t_{tab} with 95% confidence = 2.086

$$[\because \alpha = 0.05 \Rightarrow \frac{\alpha}{2} = 0.025]$$

$$t_{0.025 \text{ with } 20 \text{ d.f.}} = 2.086]$$

$$\therefore (2.086) = \frac{\bar{x} - \bar{y}}{2.72 \sqrt{\frac{1}{10} + \frac{1}{12}}}$$

$$\Rightarrow \bar{x} - \bar{y} = (2.086)(2.72) \sqrt{\frac{6+5}{60}}$$

$$\bar{x} - \bar{y} = 2.43.$$

\therefore Maximum difference between the means of samples = 2.43

7. The Blood pressure of 5 woman before and after intake of a certain drug are given below.

Before	110	120	123	132	125
After	120	118	125	136	121

Test whether there is significant change in blood pressure at 1% level of significance.

59. Let Null Hypothesis $H_0: \mu_1 = \mu_2$ 16

Alternative Hypothesis $H_1: \mu_1 < \mu_2$

Level of Significance: $\alpha = 0.01$ with 4 d.o.f ($n-1$)

Test Statistic.

$$t = \frac{\bar{d}}{s/\sqrt{n}} \quad \text{where } \bar{d} = \frac{\sum d}{n} \quad \text{and } s^2 = \frac{1}{n-1} \sum_{i=1}^n (d - \bar{d})^2 \\ = \frac{-12}{5} = -2.4$$

x (BP before intake of drug)	y (BP after intake)	d $= x - y$	$\frac{d - \bar{d}}{s}$	$(d - \bar{d})^2$
110	120	-10	-7.6	57.76
120	118	2	4.4	19.36
123	125	-2	0.4	0.16
132	136	-4	-1.6	2.56
125	121	4	6.4	40.96
		$\sum d = -12$		$\sum (d - \bar{d})^2 = 120.8$

$$\therefore s^2 = \frac{1}{4} (120.8) = 30.2$$

$$\Rightarrow s = \sqrt{30.2} = 5.4954$$

$$\therefore t = \frac{-2.4}{5.4954/\sqrt{5}} = \frac{-2.4 \sqrt{5}}{5.4954} = \frac{-5.3666}{5.4954} = -0.9765$$

But t_{tab} at $\alpha = 0.01$ with $d.o.f = 4$ is 4.604

Since $|t_{cal}| < t_{tab}$, we accept H_0

\therefore There is no significant change in BP after intake of certain Drug.

- o The average losses of workers, before and after certain program are given below. Use 0.05 level to test whether the program is effective

Before	40	70	45	120	35	55	77
After	35	65	42	116	33	50	73

H₀: Null hypothesis $H_0: \mu_1 = \mu_2$

Alternative hypothesis $H_1: \mu_1 \neq \mu_2$

level of significance: $\alpha = 0.05$ with '6' d.o.f

Test statistic

$$t_{\text{cal}} = \frac{\bar{d}}{s/\sqrt{n}} \quad \text{where } \bar{d} = \frac{\sum d}{n} = \frac{28}{7} = 4$$

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

Before (x)	After (y)	$d = x - y$	$d - \bar{d}$	$(d - \bar{d})^2$
40	35	5	1	1
70	65	5	1	1
45	42	3	-1	1
120	116	4	0	0
35	33	2	-2	4
55	50	5	1	1
77	73	4	0	0
		28		8

$$\therefore s^2 = \frac{1}{6} (8) = \frac{4}{3} = 1.333$$

$$s = \sqrt{1.333} = 1.1545$$

$$\therefore t_{\text{cal}} = \frac{4}{1.1545/\sqrt{7}} = \frac{4\sqrt{7}}{1.1545} = 9.167$$

$$\text{Here } t_{\text{tab}} = 2.447$$

$\therefore t_{\text{cal}} > t_{\text{tab}}$; we reject H_0

\therefore The program is not effective

9. Memory capacity of 10 students were tested before and after training. state whether the training was effective or not from the following scores.

17 (5)

Before	12	14	11	8	7	10	3	0	5	6
After	15	16	10	7	5	12	10	2	3	8

99. Solve as previous problems.

~~Hint: t_{tab} with $\alpha=0.05$, d.o.f = $10-1=9$ is 1.833~~

$$\begin{aligned} H_0: \mu_1 &= \mu_2 & t_{cal} &= -1.365 & \text{Accept } H_0 \\ H_1: \mu_1 &< \mu_2 & |t_{cal}| &\leq t_{tab} & \Rightarrow \text{Reject } H_0 \end{aligned}$$

∴ Training was not effective.

⇒ There is no significant difference b/w before and after training in scores.

10. The Nicotine contents in milligrams in two samples of tobacco were found to be as follows

Sample A	24	27	26	21	25	-
Sample B	27	30	28	31	22	36

So. To test whether the two samples come from same population, we have to test

- the equality of means by using Student's 't' test
- the equality of Variances by using F-test

~~Given $n_1=5$; $n_2=6$~~

Mean and S.D's of the Samples

x	$\frac{x - \bar{x}}{\sigma_x = 24.6}$	$(x - \bar{x})^2$	y	$\frac{y - \bar{y}}{\sigma_y = 29}$	$(y - \bar{y})^2$
24	0.6	0.36	27	-2	4
27	2.4	5.76	30	1	1
26	1.4	1.96	28	-1	1
21	3.6	12.96	31	2	4
25	0.4	0.16	22	-7	49
\bar{x}		21.2	$\bar{y} = 29$		108
128					

$$\therefore \bar{x} = \frac{\sum x}{n_1} = \frac{128}{5} = 24.6$$

$$\bar{y} = \frac{\sum y}{n_2} = \frac{174}{6} = 29$$

$$\sum (x - \bar{x})^2 = 21.2 \quad \text{and} \quad \sum (y - \bar{y})^2 = 108$$

$$S_1^2 = \frac{1}{n_1 - 1} \sum (x - \bar{x})^2 \\ = \frac{21.2}{4} = 5.3$$

$$S_2^2 = \frac{1}{n_2 - 1} \sum (y - \bar{y})^2 \\ = \frac{108}{5} = 21.6$$

(i) F-Test:

Null Hypothesis $H_0: \sigma_1^2 = \sigma_2^2$

Alternative Hyp $H_1: \sigma_1^2 \neq \sigma_2^2$

Level of Significance: $\alpha = 0.05$ with (5, 4) d.o.f

Test Statistic

$$F_{\text{cal}} = \frac{S_2^2}{S_1^2} \quad (\because S_2 > S_1) \\ = \frac{21.6}{5.3} = 4.075$$

But $F_{\text{tab}} = 6.26$; $\therefore F_{\text{cal}} < F_{\text{tab}}$, we accept H_0

\therefore The population (normal) have equal variances

(ii) Student's t-test.

Null hyp $H_0: \mu_1 = \mu_2$

Alternative hyp $H_1: \mu_1 \neq \mu_2$

level of significance: $\alpha = 0.05$ with 9 d.o.f

Test statistic

$$\begin{aligned} F_{\text{test}}^2 &= \frac{1}{n_1 + n_2 - 2} [\sum (m_i - \bar{m})^2 + \sum (y_j - \bar{y})^2] \\ &= \frac{1}{9} [21.2 + 108] \\ &= \frac{129.2}{9} = 14.35 \end{aligned}$$

$$\therefore f = \sqrt{14.35} = 3.78$$

$$\text{Now } t_{\text{cal}} = \frac{\bar{x} - \bar{y}}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} = \frac{24.6 - 28}{3.78 \sqrt{\frac{1}{5} + \frac{1}{6}}} = -1.92$$

$$\text{But } t_{\text{tab}} = 2.262$$

$\because |t_{\text{cal}}| < t_{\text{tab}}$ we accept H_0

\therefore The normal population come from equal mean.

- II. Pumpkins were grown under two experimental condys
 Two random samples of 11 and 10, show the sample
 S.D of their weights as 0.8 and 0.5 resp. Assume
 that the weight distribution are normal, test the
 hypothesis is that the true variances are equal.

Given $n_1 = 11$ and $n_2 = 10$

$$S_1 = 0.8 \quad S_2 = 0.5$$

$$\Rightarrow s_1^2 = \frac{1}{n_1 - 1} \sum (x_i - \bar{x})^2 = \frac{n_1 S_1^2}{n_1 - 1} = \frac{11(0.8)^2}{10} \Rightarrow s_1^2 = \sqrt{\frac{11(0.8)^2}{10}} \\ = 0.704 \quad = 0.830$$

$$S_2^2 = \frac{1}{n_2-1} \sum (y_i - \bar{y})^2 = \frac{n_2 s_1^2}{n_2-1} = \frac{10}{9} (0.5)^2 = 0.5555$$

$$\Rightarrow S_2^2 = \sqrt{\frac{10}{9}} (0.5) = 0.5555 \approx 0.5555$$

$\therefore S_2^2 > S_1^2$ [Hence Null Hypothesis: $\sigma_1^2 = \sigma_2^2$
Alternative Hypothesis: $\sigma_1^2 \neq \sigma_2^2$

Test statistic Level of significance: $\alpha = 0.05$
 $F_{\text{cal}} = \frac{s_1^2}{s_2^2} = \frac{0.2777}{0.5555} = 2.535$ with (n_1-1, n_2-1)
d.f.]

But $F_{\text{tab}} = 3.14$

$\because F_{\text{cal}} < F_{\text{tab}}$, we accept H_0

\therefore True Variances are equal.

12 Four coins were tossed 160 times and the following results were obtained-

No of heads	0	1	2	3	4
observed freq	17	52	54	31	6

Under the assumption that coins are balanced, find the expected frequencies if 0, 1, 2, 3, 4 heads and test the goodness of fit.

Sol: Here probability of getting a head in a single throw = $p = \frac{1}{2} \Rightarrow q = 1 - \frac{1}{2} = \frac{1}{2}$ Prob of success = getting head
Probability of getting 0 heads = ${}^4C_0 \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^4 = \frac{1}{16}$
 \therefore Expected frequency for 0 heads
 $= E(0) = \frac{1}{16} \times 160 = 10$ [$\because N = 160$ times]
 $[P(X=r) = {}^nC_r p^r q^{n-r}]$

13. The number of automobile accidents per week in a certain community are as follows

12, 8, 20, 2, 14, 10, 15, 6, 9, 4.

Are these frequencies in agreement with the belief that accident cond^y were the same during 10 week period.

SOP let the Null Hyp H₀: The accident cond^y were same during 10 week period

Alternative Hyp H₁: Accident cond^y were not same
level of significance? $\alpha = 0.05$ with 9 d.o.f.

O _i	E _i	O _i - E _i	(O _i - E _i) ²	(O _i - E _i) ² /E _i
12	10	2	4	0.4
8	10	-2	4	0.4
20	10	10	100	10
2	10	-8	64	6.4
14	10	4	16	1.6
10	10	0	0	0
15	10	5	25	2.5
6	10	-4	16	1.6
9	10	-1	1	0.1
4	10	-6	36	3.6
100	100		26.6	

Test statistic

$$\chi^2_{\text{cal}} = \sum \frac{(O_i - E_i)^2}{E_i} = 26.6$$

if $\chi^2_{\text{tab}} = 16.916$, we reject Null Hypothesis
 \therefore the accident cond^y were not same during
 the 10 week period.

$$E(1) = N \cdot 4 \cdot \left(\frac{1}{2}\right)^3 \cdot \left(\frac{1}{2}\right)^3 = N \cdot 4 \cdot \frac{1}{16} = \frac{1}{4} \cdot 160^{\textcircled{1}} = 40$$

$$E(2) = N \cdot 4 \cdot \left(\frac{1}{2}\right)^4 = \frac{160^{\textcircled{1}}}{16} \cdot \frac{4}{256} = \frac{10 \times 4 \times 3}{256} = 60$$

$$E(3) = N \cdot 4 \cdot \left(\frac{1}{2}\right)^4 = \frac{10 \times 4!}{3!} = 10 \times 4 = 40$$

$$E(4) = N \cdot 4 \cdot \left(\frac{1}{2}\right)^4 = 160 \cdot \frac{1}{16} = 10$$

\therefore the table is

O_i	E_i	$O_i - E_i$	$(O_i - E_i)^2$	$\frac{(O_i - E_i)^2}{E_i}$
17	10	7	49	4.9
52	40	12	36	3.6
54	60	-6	36	0.6
31	40	-9	81	2.025
6	10	-4	16	1.6
160	160			12.725
$\sum O_i$	$\sum E_i$			

Let Null Hypothesis H_0 : The data follows Binomial distribution

Let Alternative Hypothesis H_1 : The data not follows B.D

Level of significance: $\alpha = 0.05$ with 4 d.o.f

[One d.o.f is lost for parameter μ]

Test statistic

$$\chi_{\text{cal}}^2 = \sum \frac{(O_i - E_i)^2}{E_i} = 12.725$$

$\because \chi_{\text{tab}}^2$ at $\alpha = 0.05$ with 4 d.o.f = 9.488

$\therefore \chi_{\text{cal}}^2 > \chi_{\text{tab}}^2$ [∴ $\chi_{\text{cal}}^2 > \chi_{\text{tab}}^2$]

\therefore we reject H_0

\therefore The data not follows Binomial Distribution

UNIT-V

Queueing Theory

Queue:- The group of items waiting for the service including the customer, who is receiving the service is called a Queue.

The places where queue is needed are Bank, post-office, Railway station, Airport, Busstop, etc.,

Customer:

Queueing theory:-

Queueing is an activity in which every individual may experience in their daily lives. It involves two parties.

(i) Customers (ii) Server.

Customers:- The person who is waiting for the service is the customer.

Server :- The person who is servicing the customers is a server.

Describe a Queueing System

A Queueing system can be completely described by the following components

i) Input pattern (or) Arrival pattern:-

The Input describes the way in which the customers arrive and join the system. The ability of Queueing system to provide service for an arriving group of

180

customers depends not only the mean arrival rate, but also on the pattern in which they arrive.

2) Service pattern:-

The way in which the server serves is known as the service pattern. It is possible that there may be single (one) server (or) several (may) servers to serve the customers in queue. The time taken to serve a customer by the server is referred to as "service time", which is usually a random variable.

The service time distribution is taken as "negative exponential distribution".

3) Queue discipline:-

It is a rule according to which customers are selected for service when queue has been formed. The most common queue discipline is the "first come, first served" (FCFS) (or) the "first in, first out" (FIFO) rule under which the customers are serviced in the strict order of their arrivals. Other queue discipline include "last in, first out" (LIFO) - rule according to which the last arrival in the system is serviced first.

4) Queue Behaviour:- The customers generally behave in four ways

Balking :- customers may not enter the queue in view of its length. This customer behaviour is referred to as "balking".

Reneging :- This occurs when a waiting customer leaves the queue due to impatience (or) importance.

Priorities :- Some customers may be served according to the priority without waiting in the queue.

Jockeying :- If a customer jumps from one queue to another, if there are more than one queue, then it is Jockeying.

Operating characteristics of a queuing system.

Some of the operational characteristics of a queuing system, that are of general interest for the evaluation of the performance of an existing queuing system and to design a new system are as follows.

i) Expected no. of customers in the system :-

It is denoted by \bar{L}_s (or) \bar{L} is average no. of customers in the system, both waiting and in service. Here $n = \text{no. of customers in the queuing system}$.

2) Expected no. of customers in the Queue :-

It is denoted by E_m (or) L_q is the Average no. of customers waiting in the Queue. Here $m = n - 1$ i.e. excluding the customer being served.

3) Expected Waiting time in the system :-

It is denoted by E_v (or) w_s is the Average total time spent by a customer in the System. It is generally taken to be the waiting time plus servicing time.

i) Expected waiting time in Queue :-

It is denoted by E_w (or) w_q is the Average time spent by a customer in the Queue before the commencement of his service.

ii) The Server utilization factor (or) busy period :-

It is denoted by $\rho (= \frac{\lambda}{\mu})$ is the proportion of

time that a server actually spends with the customer.

Here λ = the Average no. of customers arriving per unit of time and μ = the Average no. of customers completing service per unit of time.

The server utilization factor is also known as "traffic intensity (or) the clearing ratio."

Definitions :-

Queue length :-

The no. of customers waiting in the queue at any time is called the Queue length.

Average Queue length :-

No. of customers in the queue per unit time.

Waiting time :-

The time a customer (unit) has to wait in the queue till he (it) will be taken into service.

Busy period :-

The time when the server is busy in serving. So this is the time from the starting of the service unit to the last unit.

Servicing time :-

The time taken to serve a unit is the servicing time.

Idle time :-

The time when the server is not servicing any customer (unit)

Mean arrival rate :-

Average no. of arrivals in a time interval of unit length.

Mean servicing rate :-

The average no. of units served in a time interval of unity.

Traffic Intensity :-

The ratio of the mean arrival rate to the mean service rate

$$\rho = \frac{\lambda}{\mu}$$

Transient state and Steady state

Transient state :-

If the operating characteristics are depends on time. Then their system said to be in transient state.

So in this state waiting time, service time the probability distribution are dependent on time, the system is in transient.

Steady State :-

If the operating characteristics are independent of time. Then the system said to be steady state.

so in this system waiting time, service time distribution of arrivals are independent of

time.

Explosive state :- If the arrival rate is greater than

the service rate. Then the system cannot attain

the steady state. The queue length increases

rapidly with time and tend to ∞ . This state is

called explosive state.

Notations & Symbols of Queueing system.

n - no. of customers (units) in the system . Including one that is served.

m - no. of customers in the queue excluding the one that is served.

$P_n(t)$ = The transient probability that there are n - customers in the system at any time "t".

P_n - The steady state probability that there are n - customers in the system at any time "t".

λ = no. of arrivals per unit time.

μ = no. of services per unit time.

ρ = traffic Intensity $= \left(\frac{\lambda}{\mu}\right)$

L_s = Expected (Average) no. of customers in the system.

L_q = Expected (Average) no. of customers in the queue.

Probability distribution in Queueing Systems.

It is assumed that customers joining the Queueing Systems arrive in a random manner and follow a Poisson distribution (or) equivalently the customer arrival times obey Exponential distributed.

The arrival and service distributions for Poisson Queues are derived. The basic assumptions (Axioms) governing this type of Queues are stated below.

Axiom 1 :- The no. of arrivals in non-overlapping intervals is statistically independent.

i.e The process has independent increments.

Axiom 2 :- The probability of more than one arrival b/w time t & $t + \Delta t$ is $O(\Delta t)$

Axiom 3 :- The probability that an arrival occurs b/w time t & $t + \Delta t$ is equal to $\lambda \Delta t + O(\Delta t)$ where λ is a constant. Δt is an incremental element and $O(\Delta t)$ represents the terms such that

$$\lim_{\Delta t \rightarrow 0} \frac{O(\Delta t)}{\Delta t} = 0$$

A process characterised by the above axioms is called a poisson process. This is equivalent to assuming that the inter arrival time is an exponentially distributed random variable.

→ If in the above λ is replaced by μ . The no. of services per unit time, we note that the service distribution also is a poisson process with inter-service time following an exponential distribution.

Arrival distribution theorem

If the arrivals are completely random, then the probability distribution of no. of arrivals in a fixed time-interval follows a poisson distribution.

Kendall's Notation

The representation of queuing model, as given below was introduced by D.G. Kendall and then by A. Lee

(a/b/c) : (d/e)

a - arrival distribution, b - Departure distribution

c - No. of Channels, d - maximum no. of customers allowed

e - Queue discipline
i.e. capacity of the system.

a & b are represented by M or M

M = Markovian (poisson arrival & exponential distribution)

Pure Birth and Death process:-

Distribution of arrivals (or) pure Birth process:-

A Queueing model in which only arrivals are counted and no departures take place are called pure birth process.

Distribution of Departures (or) pure Death process:-

A Queueing model in which only departures are counted and no other arrivals allowed are called pure death process.

Pure Birth & Death process:-

In Queueing model, we have arrivals (or births) & departures (or deaths). Such a model (or) process is called a pure Birth & death model (or) process.

Classification of Queueing models

The Queueing models can be categorized as:-

1) Deterministic Queueing model :-

If each customer arrives at known intervals and the service time is known with certainty. The Queueing model is said to be deterministic in nature.

2) Probabilistic Queueing model :-

If the arrivals of the customers (or) the service

times of the customers (or) both of queuing system is not known with certainty & expressed only in probabilistic nature, then the queuing model is said to be probabilistic nature.

i) Model - I

(M/M/1) : (∞ /FCFS).

The first M - denotes that the arrivals are poisson.

The second M - denotes that the departures are poisson.

1 - denotes that single service channel.

∞ - denotes that the arrivals are from an infinite population and there is no limit on the system capacity.

FIFO : first in, first out.

Characteristics of (M/M/1) : (∞ /FIFO) System

1. λ = no. of arrivals per unit time.
2. μ = no. of services per unit time.
3. $\rho = \text{traffic intensity} = \frac{\lambda}{\mu}$
4. The probability that the system is busy = ρ
5. $P_n = \text{probability that there } n \text{ people in the system}$
 $= f^n(1-f)$

6. P_0 = The probability that the system is idle = $1 - \rho$
7. $P(n \geq k)$ (There are k (or) more customers in the system)
 $= P(n > k) = \rho^k$
8. The probability that the no. of customers in the system exceeds k (or) greater than k = $P(n > k) = \rho^k$
9. Average no. of customers in the system (or) Average length of the System = $L_S = \frac{\rho}{1 - \rho}$
10. Average no. of customers in the Queue (or) Average length of the Queue = $L_Q = \frac{\rho^2}{1 - \rho}$
11. Average no. of customers in the non-empty Queue
 $= E(m|m > 0) = \frac{u}{u - \lambda}$
12. Variance of Queue Length = $\frac{\rho}{(1 - \rho)^2}$
13. The probability density function (cumulative probability distribution) of waiting time distribution (excluding service time) = $\psi(w) = \lambda (1 - \frac{\lambda}{u}) e^{-(u-\lambda)w}$ if $w > 0$
14. The probability that a customer has to wait for more than 1 minute (probability that waiting time exceeds t)
 $= \int_{t}^{\infty} \frac{\lambda}{u} (u - \lambda) \cdot e^{-(u-\lambda)t} dt = \frac{1}{u} \cdot e^{-(u-\lambda)t}$

15. Average waiting time of a customer in the Queue = w_q

$$w_q = \frac{1}{u-\lambda}$$

16. Average waiting time of a customer in a non-empty queue = $E(w|w>0) = \frac{1}{u-\lambda}$

17. Average time that a customer spends in the system

$$w_s = \frac{1}{u-\lambda}$$

18) Little's formulae

$$L_s = \lambda w_s, \quad L_q = \lambda w_q$$

$$L_q = \rho L_s, \quad w_q = \rho w_s$$

STOCHASTIC PROCESSES

* Definition :- A stochastic process is a set of random variable $\{x_i\}$ (or) $\{x_t\}$ depending on some real parameters like time 't'.

* States :- The values assumed by the random variable $x(t)$ are called states.

* State Space :-

The set of all possible values of any individual member of the random process is called "state space". It is denoted by I (or) S.

→ The state space is said to be discrete if it contains a finite (or) countable infinity of points otherwise it is called continuous.

→ If the parameter set T is discrete the random process will be denoted by $\{x_n\}$ (or) $\{x(n)\}$.

→ If the parameter set T is continuous the process will be denoted by $\{x(t)\}$ (or) $\{x_t\}$.

Examples of Stochastic process

- 1) A Queuing System
- 2) Turbulent Fluid Flow
- 3) Movement of mole cules of a gas (or) liquid.

- 4) A Random Walk model.
- 5) Communication process
- 6) Gamblers Ruin problem.

Classification of Stochastic processes

Stochastic processes mainly classified into four types.
Based on time 't' and the random variable X . They are:-

1) Continuous Stochastic processes:-

If both the random variable X and time t are continuous, the stochastic process is called a "continuous Stochastic process".

2) Discrete Stochastic processes:-

If the random variable X is discrete & time 't' is continuous, the stochastic process is called a "discrete Stochastic process".

3) Discrete Stochastic Sequence:-

If both random variable X & time 't' are discrete. Then the stochastic process is called a "Discrete Stochastic Sequence".

4) Continuous Stochastic Sequence:-

If the random variable X is continuous and time 't' is discrete then the stochastic process is called "continuous

random sequences".

Deterministic stochastic process.

A random process is called a Deterministic stochastic process if future values of any sample function can be predicted from its past observations (or) past values.

Non-Deterministic stochastic process.

A Stochastic process is called non-Deterministic Stochastic process if future values of any sample function cannot be predicted from past observations (or) past values.

Stochastic process with independent increments.

If, for all $t_1 < t_2 < \dots < t_n$, the random variables $x(t_2) - x(t_1), x(t_3) - x(t_2), \dots, x(t_n) - x(t_{n-1})$ are independent. Then the process $\{x(t)\}$ is said to be a stochastic process with independent increments.

Methods of description of random process

Consider a random process $x(t)$. For any given time t_1 , the distribution associated with a random process x_1 is given by $F(x_1, t_1) = P(x(t_1) \leq x_1)$. Is called First order distribution of the random variable $x(t_1)$ and $f(x_1, t_1) = \frac{\partial}{\partial x_1} F(x_1, t_1)$ is called First order probability density function of random variable.

$\{x(t_1) \leq x_1, x(t_2) \leq x_2\}$ is
 the joint probability distribution of the random variables
 $x(t_1), x(t_2)$ is called the second order distribution of the
 random variable $f(x_1, x_2, t_1, t_2) = \frac{\partial^2}{\partial x_1 \partial x_2} F(x_1, x_2, t_1, t_2)$ is
 called the second order density function of random variable $X(t)$
 If the n th order distribution $\{x(t)\}$ is the joint distribution
 $F(x_1, x_2, \dots, x_n; t_1, t_2, \dots, t_n)$ of the random variables
 $x(t_1), x(t_2), \dots, x(t_n)$

Stationary process and non-stationary process

i) Stationary process :-

If certain probability distribution (or) average
 do not depend on time 't'. Then the random process $\{x(t)\}$
 is called stationary process.

Depending on the density functions of the
 random variables of the process, there are several types
 of stationary they are.

i) First order stationary process :-

If the first order probability density function of
 random process is independent of time, then it is called
 as First order stationary process.

$$F(x_1, t_1) = F(x_1, t_1 + \delta) \text{ where } \delta \text{ is a real number.}$$

Stationary process:-

(ii) Second order :- If the second order probability density function of random process is independent of time, Then it is called as Second order stationary process.

$$f(x_1, x_2, t_1, t_2) = f(x_1, x_2; t_1 + \delta, t_2 + \delta).$$

where δ is real number

By if the Order probability density function of random process is independent of time. Then it is called as n^{th} order stationary process.

$$f(x_1, x_2, \dots, x_n, t_1, t_2, \dots, t_n) = f(x_1, \dots, x_n, t_1 + \delta, \dots, t_n + \delta).$$

where δ is real number.

Non-Stationary process:-

A process which is not stationary is said to be non-stationary (or) evolutionary process.

Time Averages :-

Consider stochastic (or) random process $x(t)$. The Average value of $x(t)$ taken over all times is called the time average of $x(t)$. It can be expressed as

$$\bar{x} = A[x(t)] = \lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T}^T x(t) dt$$

Average values of Single Random process and two or more random process.

The Average values of Single random process (or) two (or) more random process are as follows.

1) Mean:-

Mean function of random process $x(t)$ is expected value of a number of random process.

$$\text{i.e. Mean } u(t) = E(x(t))$$

2) Auto Correlation:-

The Auto correlation of the random process $x(t)$ is the expected value of the product of any two random variables $x(t_1) + x(t_2)$ of the process and it is denoted

$$R(t_1, t_2) \text{ (or) } R_x(t_1, t_2) \text{ (or) } R_{xx}(t_1, t_2)$$

$$\therefore R(t_1, t_2) = E[x(t_1) \cdot x(t_2)]$$

3) Auto Covariance:-

The Auto-Covariance of the random process $x(t)$ is denoted by $C(t_1, t_2)$ (or) $C_x(t_1, t_2)$ (or) $C_{xx}(t_1, t_2)$ and defined as the covariance of the random variables

$$x(t_1) \neq x(t_2)$$

$$\therefore c(t_1, t_2) = E[\{x(t_1) - \mu(t_1)\} \{x(t_2) - \mu(t_2)\}]$$

$$c(t_1, t_2) = R(t_1, t_2) - \mu(t_1) \cdot \mu(t_2).$$

4) Correlation Co-efficient :-

The correlation co-efficient of random process $x(t)$ is defined as

$$f(t_1, t_2) = \frac{c(t_1, t_2)}{\sqrt{c(t_1, t_1) \cdot c(t_2, t_2)}}$$

where $c(t_1, t_1)$ = Variance of $x(t_1)$ &

$c(t_2, t_2)$ = Variance of $x(t_2)$

5) Cross-correlation :-

The cross correlation of the random process $\alpha(t)$ of $x(t)$ is defined as

$$R_{xy}(t_1, t_2) = E\{x(t_1) \cdot y(t_2)\}$$

6) Cross-covariance :-

The cross covariance of the random process $x(t)$ of $y(t)$ is defined as the covariance of two random variables $x(t_1) \neq y(t_2)$.

$$\therefore c_{xy}(t_1, t_2) = R_{xy}(t_1, t_2) - \mu_x(t_1) \cdot \mu_y(t_2)$$

7) Cross - Correlation Co-efficient

The cross - correlation co-efficient of two random processes $x(t)$ & $y(t)$ is defined as

$$\rho_{xy}(t_1, t_2) = \frac{c_{xy}(t_1, t_2)}{\sqrt{c_{xx}(t_1, t_1) \cdot c_{yy}(t_2, t_2)}}$$

$$\rho_{xy}(t_1, t_2) = \frac{c_{xy}(t_1, t_2)}{\sqrt{c_{xx}(t_1, t_1) \cdot c_{yy}(t_2, t_2)}}$$

* If $P=q$ then (i) probability of ruin $q_z = \frac{a-z}{a}$

(ii) Expected duration of the game $d_z = z(a-z)$

* If $P \neq q$ Then (i) probability of ruin $q_z = \frac{\left(\frac{q}{P}\right)^a - \left(\frac{q}{P}\right)^z}{\left(\frac{q}{P}\right)^a - 1}$

(ii) Expected duration of the game

$$d_z = \left(\frac{-a}{q-P} \right) \frac{\left(\frac{q}{P}\right)^z - 1}{\left(\frac{q}{P}\right)^a - 1} + \frac{z}{q-P}$$

i) Calculate the probability of ruin and expected duration of the game, where

i) $a=50, z=40, P=0.5$

ii) $a=100, z=5, P=0.6$

Sol (i) Given that $a=50, z=40, P=0.5 = \frac{1}{2}$

$$\begin{aligned} \therefore q &= 1-0.5 \\ &= 0.5 \\ &= \frac{1}{2} \end{aligned}$$

$$\therefore \text{probability of ruin } q_z = \frac{a-z}{a} = \frac{50-40}{50} = \frac{10}{50} = 0.2$$

$$\begin{aligned} \text{Expected duration of game } d_z &= z(a-z) \\ &= 40(50-40) \\ &= 40 \times 100 = 400 \end{aligned}$$

(ii) Given that $a = 100$; $z = 5$, $P = 0.6$

$$Q = 1 - 0.6$$

$$= 0.4$$

$$\therefore \frac{P}{Q} = \frac{2}{3}$$

$$\text{probability of win} = q_V = \frac{\left(\frac{q}{P}\right)^a - \left(\frac{q}{P}\right)^z}{\left(\frac{q}{P}\right)^a - 1}$$

$$q_V = \frac{\left(\frac{2}{3}\right)^{100} - \left(\frac{2}{3}\right)^5}{\left(\frac{2}{3}\right)^{100} - 1} = 0.132$$

expected duration of the game

$$d_z = \left(\frac{-a}{q_V - P} \right) \frac{\left(\frac{q}{P}\right)^z - 1}{\left(\frac{q}{P}\right)^a - 1} + \frac{z}{q_V - P}$$

$$= \frac{-100}{-0.2} \left(\frac{\left(\frac{2}{3}\right)^5 - 1}{\left(\frac{2}{3}\right)^{100} - 1} \right) + \frac{5}{-0.2} = 409$$

* Markov process:-

A stochastic process $\{x(t), t \geq 0\}$ is called a markov process if $P[x(t_{n+1}) \leq x_{n+1} / X(t_n) = x_n, X(t_{n-1}) = x_{n-1}, \dots, X(t_0) = x_0] = P[x(t_{n+1}) \leq x_{n+1} / X(t_n) = x_n]$.

Markov chain

A stochastic process $\{x_n : n=1, 2, \dots\}$ is called markov chain, if for $j, k, j_1, j_2, \dots, j_{n-1} \in N$ (or) any subset of I_j .

$$P[x_n = k | x_{n-1} = j, x_{n-2} = j_1, \dots, x_0 = j_{n-1}]$$

$$P[x_n = k | x_{n-1} = j] = P_{jk}$$

where j_1, j_2, \dots are called the states of the markov chain.

→ If the Transition probability P_{jk} is independent of n ;

The markov chain is called Homogenous markov chain.

→ If the Transition probability P_{jk} is dependent on n ;

The chain is said to be Non-Homogenous markov chain.

The Transition probability P_{jk} refers to the states (j, k) at two successive trials. The transition is one-step and P_{jk} is called one step (or) unit step transition probability.

* Multiple ^{steps} Transition probability:

In the more general case, we are concerned with the pair of states (j, k) at two non-successive trials, say, state j at the n^{th} trial and state k at the $(n+m)^{\text{th}}$ trial. The corresponding transition probability is called m -step transition probability and is denoted by $P_{jk}^{(m)}$.

$$\text{i.e } P_{jk}^{(m)} = P[x_{n+m} = k | x_n = j]$$

* Probability vector :-

A probability vector is a vector (a row or column matrix) which is non-negative and all elements adding upto unity.

If P_1, P_2, \dots, P_n be the set of 'n' probabilities of a variable, Then the probability vector will be

$$P = [P_1, P_2, \dots, P_n] \text{ (or)} P = \begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_n \end{bmatrix} \text{ with } \sum_{i=1}^n P_i = 1 \text{ and } 0 < P_i < 1$$

$$\forall i = 1, 2, \dots, n.$$

Transition matrix (or) Matrix of Transition probabilities.

The Transition probabilities p_{jk} satisfy

$$(i) p_{jk} \geq 0$$

$$(ii) \sum_k p_{jk} = 1 \quad \forall j$$

These probabilities may be written in the matrix form

$$P = \begin{bmatrix} P_{11} & P_{12} & P_{13} & \cdots \\ P_{21} & P_{22} & P_{23} & \cdots \\ P_{31} & P_{32} & P_{33} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

This is called the transition probability matrix (or) matrix of transition probabilities (t.p.m) of the markov-chain.

* Stochastic Matrix :-

A stochastic matrix 'P' is a square matrix with non-negative elements and unit row sums.

* If P & Q are Stochastic matrices then product PQ is also a stochastic matrix. Thus P^n is a stochastic matrix for all positive integer values of n . ($\forall n \in \mathbb{Z}^+$)

Order of a Markov chain :-

A markov chain $\{x_n\}$ is said to be of order s ($s=1, 2, \dots$)

If for all n ,

$$P[x_n = k | x_{n-1} = i_1, x_{n-2} = i_2, \dots, x_{n-s} = i_{s-1}]$$

$$= P(x_n = k | x_{n-1} = i_1, \dots, x_{n-s} = i_{s-1})$$

When ever the L.H.S is defined.

→ A markov chain $\{x_n\}$ is said to be of order - one

(or simply a markov chain) if

$$P[x_n = k | x_{n-1} = i_1, x_{n-2} = i_2, \dots]$$

$$= P[x_n = k | x_{n-1} = i]$$

$$= P_{jk}$$

When ever $P[x_{n-1} = i_1, x_{n-2} = i_2, \dots] \neq 0$

→ A chain is said to be of order zero if $P_{jk} = P_k \forall j$

This implies independence of x_n & x_{n-1} .

Finite markov chain :-

A markov chain $\{x_n, n \geq 0\}$ with K states, where K is finite, is said to be a finite markov chain.

The transition matrix p is, in this case, a square with K rows and K columns.

Ex :- 1) A particle performs a random walk with absorbing barriers at 0 & 4. When ever it is at any position r ($0 < r < 4$), it moves to $r+1$ with probability p or to $(r-1)$ with probability q if $p+q=1$. But as soon as it reaches 0 or 4, it remains there itself. Let x_n be the position of the particle after n moves. The different states of x_n are the different positions of the particle. $\{x_n\}$ is a markov chain whose unit-step transition probabilities are given by:

$$P[x_n = r+1 / x_{n-1} = r] = p$$

$$P[x_n = r-1 / x_{n-1} = r] = q \quad 0 < r < 4$$

$$P[x_n = 0 / x_{n-1} = 0] = 1$$

$$P[x_n = 4 / x_{n-1} = 4] = 1$$

The Transition matrix is given by

$$\begin{array}{c} \text{States of } x_n \\ \begin{array}{ccccc} 0 & 1 & 2 & 3 & 4 \end{array} \\ \left[\begin{array}{ccccc} 1 & 0 & 0 & 0 & 0 \\ q & 0 & p & 0 & 0 \\ 0 & q & 0 & p & 0 \\ 0 & 0 & q & 0 & p \\ 0 & 0 & 0 & 0 & 1 \end{array} \right] \\ \text{States of } x_{n-1} \end{array}$$

Denumerably infinite (or) Countably infinite markov chain.

The no. of states could however be infinite when the possible values of x_n form denumerable set, Then the markov chain to be denumerably infinite (or) Countably infinite and the chain is said to have a countable state space.

Ex:- Suppose that a coin with probability p for a head is tossed indefinitely. Let x_n ; The outcome of the n^{th} trial be k : where $k (=0, 1, \dots)$ denotes that there is a run of k successes. i.e the length of the uninterrupted block of heads is k . $\{x_n, n \geq 0\}$ constitutes a markov chain with unit step transition probabilities

$$P_{jk} = P[x_n = k \mid x_{n-1} = j] = \begin{cases} p, & k = j+1 \\ q, & k = 0 \\ 0, & \text{otherwise} \end{cases}$$

The Transition matrix is given by

States of X_n

$$\begin{matrix} & 0 & 1 & 2 & \dots & K & X_{n+1} \\ 0 & q_1 & p & 0 & \dots & 0 & \dots \\ 1 & q_2 & 0 & p & \dots & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ k & q_k & 0 & 0 & \dots & 0 & p \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \end{matrix}$$

States of X_{n+1}

Markov chains as Graphs

The states of markov chain may be represented by the vertices (nodes) of the graph of one-step transitions b/w states by directed arcs: If $i \rightarrow j$, then the two vertices i & j are joined by directed arc with arrow towards j . The value of P_{ij}^n , which corresponds to the arc weight, may be indicated in the directed arc. If $S = \{1, 2, \dots, m\}$ is the set vertices corresponding to the state space of the chains 'a' is the set of directed arcs b/w these vertices;

Then the graph $G = \{S, a\}$ is directed graph (or) digraph (or) Transition graph of the chain. A digraph such that its arc weights are positive & sum of the arc

Weights of the arcs from each node unity is called a Stochastic graph; The diagram of a markov chain is a stochastic graph.

- A path in a diagraph is any sequence of arcs where the final vertex of one arc is the initial vertex of the next arc.

Higher Transition probabilities.

- * Chapman - Kolmogorov theorem.

If P is the tpm of a homogeneous markov chain,
Then the n -step tpm $p(n)$ is equal to P^n

$$\text{i.e. } [P_{ij}^{(n)}] = [P_{ij}]^n$$

- * Classification of states & chains

The states j , $j=0, 1, 2, \dots$ of a markov chain $\{x_n, n \geq 0\}$ can be classified in a distinctive manner according to some fundamental properties of the system.

Communication relation

- i) If $P_{ij}^{(n)} > 0$ for some $n \geq 1$, then we say that state j is accessible from state i , the relation is denoted

by $i \rightarrow j$ conversely if $\forall n$; $p_{ij}^{(n)} = 0$, then j is not accessible from i , in notation $i \not\rightarrow j$

2) Two accessible states i & j are said to be communicate state i communicates with itself $\forall i \geq 0$. If the state i communicates with state j & state j communicates with state k ; Then state i communicates with state k . Two states that communicate are in the same class. A state is called an essential state. If it communicates with every state it leads so.

3) If $p_{ij}^{(n)} > 0$ for some n . $\forall i \neq j$. Then every state can be reached from every other state. Then the markov chain is said to be irreducible. Then the transition matrix is irreducible. otherwise the chain is said to be reducible (or) non-irreducible.

4) A state i is said to be an absorbing state. If and only if $p_{ii} = 1$. A markov chain is absorbing if it has at least one absorbing state & it is possible to go from every non-absorbing state to atleast one absorbing state in one (or) more steps.

5) periodicity :- A state i is a returnstate if $p_{ii}^{(n)} > 0$ for some $n \geq 1$. The period d_i of a return to state i is defined as the greatest common divisor of all m such that $p_{ii}^{(m)} > 0$. Thus

$$d_i = \text{G.C.D} \{ m : p_{ii}^{(m)} > 0 \}$$

A State i is said to aperiodic if $d_i = 1$ & periodic if $d_i \geq 1$. Clearly state i is aperiodic if $p_{ii} \neq 0$.

Classification of States

7) The probability that the chain starting from state i returns to state i for the first time at the n^{th} step (or) after n (transitions) is denoted by $f_{ii}^{(n)}$, $n=1, 2, \dots$ and is called as First time return time probability (or) the recurrence time probability.

If $F_{ii} = \sum_{n=1}^{\infty} f_{ii}^{(n)} = 1$, the return to state i is certain & if $M_{ii} = \sum_{n=1}^{\infty} n f_{ii}^{(n)}$ is called the mean recurrence time of the state i .

Persistent (or) recurrent state :-

A State i is said to be persistent (or) recurrent if $F_{ii} = 1$.

Transient State :-

The state "p" is said to be transient if $\pi_{pp} < 1$
(i.e. the return to state p is uncertain)

Null persistent & Non-Null persistent State :-

The state "p" is said to be non-null persistent (or) positive persistent if π_{pp} is finite. The state "p" is said to be Null persistent if $\pi_{pp} = \infty$.

Classification of Chains

8) Ergodic :-

A positive recurrent and aperiodic state of a markov chain is called ergodic. A markov chain all of whose states are ergodic said to be a ergodic chain.

9) AtPM, P is said to be a regular matrix if all the entries of power P^m ($m = 2, 3, \dots$) are positive. A homogenous markov chain is said to be regular chain if its tPM is a regular matrix.

AtPM, P is said to be stochastic matrix, if the elements of each of the rows are non-negative and the sum of elements in each row is equal to 1.

A homogeneous markov chain will have a tpm. That is independent of initial state i after steps n as $n \rightarrow \infty$ and is called steady state probability.

$$\lim_{n \rightarrow \infty} p_{ij}^{(n)} = \pi_j \quad \forall j$$

called limiting state probability and is interpreted as the long-run proportion of time the markov chain spends in state j .

Theorem :- *

A stochastic matrix P is not regular if it occurs in the principle main diagonal.

1. Which of the following matrices are stochastic.

(i) $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ (ii) $\begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$ (iii) $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$ (iv) $\begin{bmatrix} 1 & 0 \\ -1 & 0 \end{bmatrix}$ (v) $\begin{bmatrix} 0 & 2 \\ \frac{1}{4} & \frac{1}{4} \end{bmatrix}$

Sol (i) The matrix is a square matrix with non-negative entries and sum of elements in each row is equal to 1

$\therefore \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ & $\begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$ are stochastic matrices

(iii) It is not a square matrix

\therefore It is not stochastic.

(iv) The matrix is not stochastic, because it contains negative elements

(v) The matrix is square matrix but sum in each row is not equal to 1

\therefore It is not a stochastic matrix.

8) Which of the Stochastic matrices are regular.

(i) $A = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$ (ii) $\begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix}$ (iii) $\begin{bmatrix} 0 & 0 & 1 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 1 & 0 \end{bmatrix}$

(iv) $\begin{bmatrix} \frac{1}{2} & \frac{1}{4} & 1 \\ 0 & \frac{1}{2} & 1 \\ 0 & 0 & 1 \end{bmatrix}$

Sol (ii) $G \cdot T \cdot A = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$

Not Regular since 1 lies on the main diagonal

(ii) $G \cdot T \cdot B = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix}$

$B^2 = B \cdot B = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{3}{8} & \frac{3}{8} & 0 \end{bmatrix}$

$B^3 = B^2 \cdot B = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{7}{16} & \frac{7}{16} & \frac{1}{8} \end{bmatrix}$

Since entries b_{13}, b_{23} are zero

$\therefore B$ is not regular.

$$(iii) G.T. C = \begin{bmatrix} 0 & 0 & 1 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 1 & 0 \end{bmatrix}$$

$$\therefore C^2 = \begin{bmatrix} 0 & 1 & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}, \quad C^3 = \begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$C^4 = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \end{bmatrix}, \quad C^5 = \begin{bmatrix} \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{8} & \frac{1}{2} & \frac{3}{8} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix}$$

Since all entries of some powers of C are positive

$\therefore C$ is regular stochastic process

(iv) G.T

$$D = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} & 1 \\ 0 & \frac{1}{2} & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

Not regular since 1 lies on the main diagonal.

3) Represent the following transition matrices as a diagram

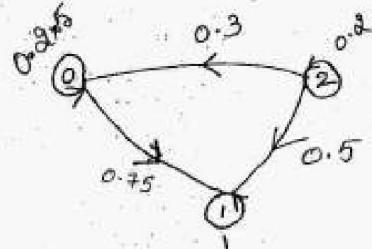
$$(i) \begin{bmatrix} 0.25 & 0.75 & 0 \\ 0 & 1 & 0 \\ 0.3 & 0.5 & 0.2 \end{bmatrix}$$

$$(ii) \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 0 & \frac{3}{4} & \frac{1}{4} \end{bmatrix}$$

Sol:- (i) G.T the tpm of markov chain is

$$P = \begin{bmatrix} 0 & 1 & 2 \\ 0 & 0.25 & 0.75 & 0 \\ 1 & 0 & 1 & 0 \\ 2 & 0.3 & 0.5 & 0.2 \end{bmatrix}$$

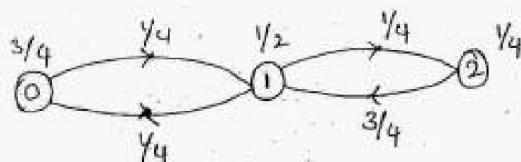
The diagram of chain is



(ii) G.T, TPM of markov chain is

$$P = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & 0 \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 0 & \frac{3}{4} & \frac{1}{4} \end{bmatrix}$$

∴ The diagram of chain is



- a) Suppose that the probability of a dry day (state 0) following a rainy day (state 1) is $\frac{1}{3}$ and that the probability of a rainy day following a dry day is $\frac{1}{2}$. Given that may 1 is a rainy day. Find the probability (i) may 3 is also a dry day
& (ii) may 5 is also a dry day.

$$\begin{array}{l} \text{dry day} = 0 \\ \text{rainy day} = 1 \end{array}$$

Sol:- we have two-state markov chain such that

$$P_{1,0} = \frac{1}{3} \text{ & } P_{0,1} = \frac{1}{2} \text{ if tpm is given by}$$

$$P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

we have $P^2 = \begin{bmatrix} \frac{5}{12} & \frac{7}{12} \\ \frac{7}{18} & \frac{11}{18} \end{bmatrix}$

$$P^4 = \begin{bmatrix} \frac{173}{432} & \frac{259}{432} \\ \frac{259}{648} & \frac{389}{648} \end{bmatrix}$$

Given that, May 1 is a dry day, then the probability that May 3 is a dry day is $\frac{5}{12}$ if May 5 is a dry day is $\frac{173}{432}$

- 5) A training process is considered as a two-state markov chain. If it rains, it is considered to be in state 0 & if it does not rain, the chain is in the state of 1. The transition probability of the markov chain is defined by $P = \begin{bmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix}$. Find the probability that it will rain for 3 days from today assuming that it is raining today. Assume that the mutual probabilities of state 0 (or) state 1 as 0.4 & 0.6 respectively.

Sol:- G.T. the one step t.p.m is given by

$$P = \begin{bmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix}$$

$$\begin{aligned} P^2 &= P \cdot P \\ &= \begin{bmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix} \begin{bmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix} \end{aligned}$$

(14)

$$P^2 = \begin{bmatrix} 0.44 & 0.56 \\ 0.38 & 0.72 \end{bmatrix}$$

$$P^3 = P^2 \cdot P = \begin{bmatrix} 0.44 & 0.56 \\ 0.38 & 0.72 \end{bmatrix} \begin{bmatrix} 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix}$$

$$P^3 = \begin{bmatrix} 0.376 & 0.624 \\ 0.312 & 0.688 \end{bmatrix}$$

The probability that it will rain on 3rd day, given that it will rain today is 0.376.

6) Let $\{x_n, n \geq 0\}$ be a markov chain with

Three states $\{0, 1, 2\}$ & with transition matrix

$$P = \begin{bmatrix} 0 & 1 & 2 \\ 0 & \frac{3}{4} & \frac{1}{4} & 0 \\ 1 & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 2 & 0 & \frac{3}{4} & \frac{1}{4} \end{bmatrix} \text{ & the initial distribution}$$

$$P\{x_0 = i\} = \frac{1}{3}, i = 0, 1, 2 \quad \text{Find } P\{x_3 = 1, x_2 = 2, x_1 = 1, x_0 = 2\}$$

Sol: - G.T, the TPM of the Markov Chain is

$$P = \begin{bmatrix} 0 & 1 & 2 \\ 0 & \frac{3}{4} & \frac{1}{4} & 0 \\ 1 & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ 2 & 0 & \frac{3}{4} & \frac{1}{4} \end{bmatrix} \text{ & } P\{x_0 = i\} = \frac{1}{3}, i = 0, 1, 2$$

Then we have

$$P\{x_1 = 1 / x_0 = 2\} = \frac{1}{4} = p_{21}$$

$$P\{x_2 = 2 / x_1 = 1\} = \frac{1}{4} = p_{12}$$

$$P\{x_2=2, x_1=1/x_0=2\} = P\{x_2=2/x=1\} \cdot P\{x_1=1/x_0=2\}$$

$$= \frac{1}{4} \cdot \frac{3}{4} = \frac{3}{16}$$

$$\therefore P\{x_2=2, x_1=1, x_0=2\} = P\{x_2=2/x_0=2\} \cdot P\{x_0=2\}$$

$$= \cancel{\frac{1}{16}} \cdot \cancel{\frac{1}{16}} = \frac{1}{16}$$

$$\therefore P\{x_3=1, x_2=2, x_1=1, x_0=2\} = P\{x_3=1/x_2=2, x_1=1, x_0=2\}$$

$$= P\{x_3=1/x_2=2\} \cdot P\{x_2=2, x_1=1, x_0=2\}$$

$$= P\{x_3=1/x_2=2\} \cdot P\{x_2=2, x_1=1, x_0=2\}$$

$$= \frac{3}{4} \cdot \frac{1}{16} = \frac{3}{64}$$

Q) The t.p.m of a markov chain $\{x_n\}$, $n=1, 2, 3, \dots$ having

3 states 1, 2 & 3 is $P = \begin{bmatrix} 0.1 & 0.5 & 0.4 \\ 0.6 & 0.2 & 0.2 \\ 0.3 & 0.4 & 0.3 \end{bmatrix}$ and the initial

distribution is $p^{(0)} = (0.7, 0.2, 0.1)$ find (i) $P\{x_2=3\}$ (ii) $P\{x_3=2, x_2=3, x_1=3, x_0=2\}$

Sol G.T the t.p.m of markov chain is

$$P = \begin{bmatrix} 1 & 2 & 3 \\ 0.1 & 0.5 & 0.4 \\ 0.6 & 0.2 & 0.2 \\ 0.3 & 0.4 & 0.3 \end{bmatrix}$$

We have $P(x_0=1)=0.7, P(x_0=2)=0.2, P(x_0=3)=0.1$

$$\text{Also } P^{(2)} = P^2 = \begin{bmatrix} 0.1 & 0.5 & 0.4 \\ 0.6 & 0.2 & 0.2 \\ 0.3 & 0.4 & 0.3 \end{bmatrix} \begin{bmatrix} 0.1 & 0.5 & 0.4 \\ 0.6 & 0.2 & 0.2 \\ 0.3 & 0.4 & 0.3 \end{bmatrix}$$

$$P^2 = \begin{matrix} & 1 & 2 & 3 \\ 1 & \begin{bmatrix} 0.43 & 0.31 & 0.26 \end{bmatrix} \\ 2 & \begin{bmatrix} 0.24 & 0.42 & 0.34 \end{bmatrix} \\ 3 & \begin{bmatrix} 0.36 & 0.35 & 0.29 \end{bmatrix} \end{matrix}$$

$$\begin{aligned}
(i) \quad P\{X_2=3\} &= \sum_{i=1}^3 P\{X_2=3/X_0=i\} \cdot P\{X_0=i\} \\
&= P\{X_2=3/X_0=1\} \cdot P\{X_0=1\} + P\{X_2=3/X_0=2\} \\
&\quad P\{X_0=2\} + P\{X_2=3/X_0=3\} \cdot P\{X_0=3\} \\
&= P_{13}^{(2)} \cdot P\{X_0=1\} + P_{23}^{(2)} \cdot P\{X_0=2\} + P_{33}^{(2)} \cdot P\{X_0=3\} \\
&= 0.26(0.7) + 0.34(0.2) + 0.29(0.1) \\
&= 0.279.
\end{aligned}$$

$$\begin{aligned}
(ii) \quad P\{X_1=3/X_0=2\} &= P_{23} = 0.2 \\
P\{X_1=3, X_0=2\} &= P\{X_1=3/X_0=2\} \cdot P\{X_0=2\} \\
&= 0.2 \times 0.2 \\
&= 0.04 \\
P\{X_2=3, X_1=3, X_0=2\} &= P\{X_2=3/X_1=3, X_0=2\} \cdot P\{X_1=3, X_0=2\} \\
&= P\{X_2=3/X_1=3\} \cdot P\{X_1=3, X_0=2\} \\
&= 0.3 \times 0.04 \quad (\text{by markov property}) \\
&= 0.012
\end{aligned}$$

$$\begin{aligned}
 & P\{x_3=2, x_2=3, x_1=3, x_0=2\} \\
 & = P\{x_3=2 / x_2=3, x_1=3, x_0=2\} \cdot P\{x_2=3, x_1=3, x_0=2\} \\
 & = P\{x_3=2 / x_2=3\} \cdot P\{x_2=3, x_1=3, x_0=2\} \\
 & = 0.4 \times 0.012 \\
 & = 0.0048
 \end{aligned}$$

- 8) A fair die tossed repeatedly. If x_n denotes the maximum of the numbers occurring in the first n tosses. Find the transition probability matrix P of the markov chain $\{x_n\}$.
 Find also p^* and $P(x_2=6)$

Sol:- State Space $= \{1, 2, 3, 4, 5, 6\}$

The TPM is formed using the following analysis

Let x_n = The maximum of the numbers occurring in the first n trials = 3 (say)

Then $x_{n+1} = 3$. If the $(n+1)^{\text{th}}$ trial results is 1, 2 or 3.

= 4. If the $(n+1)^{\text{th}}$ trial results is 4

= 5. If the $(n+1)^{\text{th}}$ trial results is 5

= 6. If the $(n+1)^{\text{th}}$ trial results is 6

$$\therefore P\{x_{n+1}=3 / x_n=3\} = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2}$$

$$\therefore P\{x_{n+1}=i / x_n=3\} = \frac{1}{6} \text{ When } i = 4, 5, 6$$

∴ The tpm of chain is

$$P = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 1 & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ 2 & 0 & \frac{2}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ 3 & 0 & 0 & \frac{5}{6} & \frac{1}{6} & \frac{1}{6} \\ 4 & 0 & 0 & 0 & \frac{9}{6} & \frac{1}{6} \\ 5 & 0 & 0 & 0 & 0 & \frac{5}{6} \\ 6 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\frac{6}{6} = 1$$

p_0
 $p^{(1)}$
 $p^{(2)}$

$$(ii) P^2 = \frac{1}{36} \begin{bmatrix} 1 & 3 & 5 & 7 & 9 & 11 \\ 0 & 4 & 5 & 7 & 9 & 11 \\ 0 & 0 & 9 & 7 & 9 & 11 \\ 0 & 0 & 0 & 16 & 9 & 11 \\ 0 & 0 & 0 & 0 & 25 & 11 \\ 0 & 0 & 0 & 0 & 0 & 36 \end{bmatrix}$$

Initial state probability distribution is $p^{(0)} = (\frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6})$

(Since all the values of 1, 2, ..., 6 are equally likely)

$$\begin{aligned}
 (ii) P(x_2=6) &= \sum_{i=1}^6 P(x_2=6/x_0=i) P(x_0=i) \\
 &= \frac{1}{6} \sum_{i=1}^6 P_{i6}^{(2)} \\
 &= \frac{1}{6} \left[P_{16}^{(2)} + P_{26}^{(2)} + P_{36}^{(2)} + P_{46}^{(2)} + P_{56}^{(2)} + P_{66}^{(2)} \right] \\
 &= \frac{1}{6} \left[\frac{11}{36} + \frac{11}{36} + \frac{11}{36} + \frac{11}{36} + \frac{11}{36} \right] \\
 &= \frac{1}{6 \times 36} [11 + 11 + 11 + 11 + 11 + 36] = \frac{1}{36} \times 56 = \frac{14}{9}
 \end{aligned}$$

$$P(X_2 = 6) = \frac{91}{816}$$

9) A gambler has Rs. 2. He bets Rs. 1 at a time and wins Rs. 1 with probability $\frac{1}{2}$. He stops playing if he loses Rs. 2.

(or) Wins Rs. 4.

(a) What is the tpm of related markov chain?

(b) What is the probability that he has lost his money at the end of 5 plays?

(c) What is the probability that the game lasts more than 7 plays?

Sol Let X_n represent the amount with the player at the end of the n^{th} round of the play.

The state space $X_n = \{0, 1, 2, 3, 4, 5, 6\}$

When the game is stopped, if the player loses Rs. 2, $X_n = 0$

(or) If he wins Rs. 4, $X_n = 6$.

The tpm of markov chain is written as

	0	1	2	3	4	5	6
0	1	0	0	0	0	0	0
1	$\frac{1}{2}$	0	$\frac{1}{2}$	0	0	0	0
2	0	$\frac{1}{2}$	0	$\frac{1}{2}$	0	0	0
3	0	0	$\frac{1}{2}$	0	$\frac{1}{2}$	0	0
4	0	0	0	$\frac{1}{2}$	0	$\frac{1}{2}$	0
5	0	0	0	0	$\frac{1}{2}$	0	$\frac{1}{2}$
6	0	0	0	0	0	0	1

(11)

This markov chain is called random walks with absorbing barriers at 0 & 6, since the chain cannot come out of the states 0 & 6, once it has entered.

b) Since the player has Rs. 2 to start the play, the initial probability distribution of x_0 is

$$p^{(0)} = (0 \ 0 \ 1 \ 0 \ 0 \ 0)$$

probability distribution after one play is given by

$$p^{(1)} = p^{(0)} \cdot P = (0 \ 0 \ 1 \ 0 \ 0 \ 0) P = (0 \ \frac{1}{2} \ 0 \ \frac{1}{2} \ 0 \ 0 \ 0)$$

$$p^{(2)} = p^{(1)} \cdot P = (0 \ \frac{1}{2} \ 0 \ \frac{1}{2} \ 0 \ 0 \ 0) P = (\frac{1}{4} \ 0 \ \frac{1}{2} \ 0 \ \frac{1}{4} \ 0 \ 0)$$

$$p^{(3)} = p^{(2)} \cdot P = (\frac{1}{4} \ 0 \ \frac{1}{2} \ 0 \ \frac{1}{4} \ 0 \ 0) P = (\frac{1}{4} \ \frac{1}{4} \ 0 \ \frac{3}{8} \ 0 \ \frac{1}{8} \ 0)$$

$$p^{(4)} = p^{(3)} \cdot P = (\frac{1}{4} \ \frac{1}{4} \ 0 \ \frac{3}{8} \ 0 \ \frac{1}{8} \ 0) P = (\frac{3}{8} \ 0 \ \frac{5}{16} \ 0 \ \frac{1}{4} \ 0 \ \frac{1}{16})$$

$$p^{(5)} = p^{(4)} \cdot P = (\frac{3}{8} \ 0 \ \frac{5}{16} \ 0 \ \frac{1}{4} \ 0 \ \frac{1}{16}) P = (\frac{3}{8} \ \frac{5}{32} \ 0 \ \frac{9}{32} \ 0 \ \frac{1}{8} \ \frac{1}{16})$$

The probability that the player has lost his money at the end of 5 plays $= P(x_5 = 0) = \frac{3}{8}$

c) The probability that the game lasts more than 7 plays

$= P(\text{system is neither in state 0 nor in 6 at the end of the seventh round})$

$$\text{using } P^{(6)} = P \cdot P = \left(\begin{array}{cccccc} \frac{29}{64} & 0 & \frac{7}{32} & 0 & \frac{13}{64} & 0 & \frac{1}{8} \end{array} \right)$$

$$P^{(7)} = P^{(6)} \cdot P = \left(\begin{array}{cccccc} 0 & 1 & \frac{23}{32} & \frac{9}{16} & \frac{13}{128} & 0 & \frac{1}{8} \\ \frac{29}{64} & \frac{1}{64} & 0 & \frac{27}{128} & 0 & \frac{13}{128} & \frac{1}{8} \end{array} \right)$$

$$\therefore P\{x_7 = 1, 2, 3, 4 \text{ or } 5\} = \frac{1}{64} + 0 + \frac{27}{128} + 0 + \frac{13}{128} \\ = \frac{27}{64}$$

- 10) Three boys A, B & C are throwing a ball to each other. A always throws the ball to B & B always throws the ball to C; but C is just as likely to throw the ball to B as A. Show that the process is markovian. Find the transition matrix & classify the states. Do all the states are ergodic.

Sol:- The transition probability matrix of the process $\{x_n\}$ is given below.

$$P_{(x_n, y_n)} = \begin{matrix} \text{State of } x_n (\text{GIVEN}) \\ \begin{array}{c} A \\ B \\ C \end{array} \end{matrix} \quad \begin{matrix} \text{State of } x_{n-1} \\ \begin{array}{c} A \\ B \\ C \end{array} \end{matrix} = \begin{matrix} A & \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \\ B & \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \\ C & \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} \end{matrix}$$

State of x_n depends only on state of x_{n-1} but not on state of x_{n-2}, x_{n-3}, \dots or earlier states.

$\therefore \{x_n\}$ is a markov chain.

$$\text{Now } P^2 = \begin{bmatrix} 0 & 0 & 1 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}, \quad P^3 = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix}$$

$P_{11}^{(3)} > 0, P_{13}^{(2)} > 0, P_{21}^{(2)} > 0, P_{22}^{(2)} > 0, P_{33}^{(2)} > 0$ and all other

∴ The chain is Irreducible.

Irreducible,

$$P^4 = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \end{bmatrix} \quad P^5 = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{8} & \frac{3}{8} & \frac{1}{2} \end{bmatrix} \text{ a.c.t.b.l.e.}$$

$$P_{ij}^{(n)} > 0.$$

$$P^6 = \begin{bmatrix} \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & \frac{3}{8} & \frac{1}{2} \\ \frac{1}{8} & \frac{3}{8} & \frac{3}{8} \end{bmatrix} \text{ & so on.}$$

We note that $P_{11}^{(2)}, P_{11}^{(3)}, P_{11}^{(4)}, P_{11}^{(5)}, P_{11}^{(6)}$ etc. are > 0 . for $i=2,3$

G.C.D of 2, 3, 4, 5, 6, ... = 1

The states 2 & 3 (i.e. B & C) are periodic with period 1
i.e. a periodic.

We note that $P_{11}^{(3)}, P_{11}^{(5)}, P_{11}^{(6)}$ etc., are > 0 . & G.C.D of
3, 5, 6, ... = 1

∴ The state 1 (i.e. state A) is periodic with period 1 i.e.
aperiodic.

Since the chain is finite and irreducible, all its

States are non-null persistent. (Q9)

Hence, all the states are ergodic.

- ii) Find the nature of states of the markov chain with
+ P.M

$$P = \begin{bmatrix} 0 & 1 & 2 \\ 0 & 1 & 0 \\ 1 & 0 & \frac{1}{2} \\ 2 & 0 & 1 \end{bmatrix}$$

Sol:- G.T the time of markov chain Ps

$$P = \begin{bmatrix} 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 1 & 0 \end{bmatrix}$$

$$P^2 = P \cdot P = \begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$$

$$P^3 = P^2 \cdot P = \begin{bmatrix} 0 & 1 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 1 & 0 \end{bmatrix} = P$$

$$P^4 = P^3 = P = P \cdot P = P^2$$

$$\text{Hence } P^5 = P^3 \cdot P^2 = P \cdot P^2 = P$$

$$P^6 = P^4 \cdot P^2 = P^2 \cdot P^2 = P^2$$

$$\text{In general } P^{2n} = P^2 \text{ & } P^{2n+1} = P^2 \cdot P = P \cdot P = P^3 = P$$

Also, we observe

$$P_{00}^{(2)} > 0, P_{02}^{(2)} > 0, P_{11}^{(2)} > 0, P_{20}^{(2)} > 0, P_{22}^{(2)} > 0.$$

$$P_{01}^{(1)} > 0, P_{10}^{(1)} > 0, P_{12}^{(1)} > 0, P_{21}^{(1)} > 0.$$

Hence the markov chain is irreducible.

Also $P_{ii}^{(2)} > 0, P_{ii}^{(4)} > 0, P_{ii}^{(6)} \geq 0$ & so on for every 'i'.

Thus $P_{ii}^{(2n)} > 0, P_{ii}^{(2n+1)} = 0$, for each i.

\therefore All the states of the markov chain are periodic with period 2.

Since the markov chain is finite & irreducible, all its states are non-null persistent. All the states are not ergodic.

* Irreducible:- If $P_{ij}^{(n)} > 0$ for some n & $\forall i \neq j$. Then every state can be reached from every other state. When this condition is satisfied, the markov chain is said to be irreducible. The TPM of an irreducible chain is an irreducible matrix otherwise the chain is said to be reducible.

19) The transition probability matrix of a markov chain is given

by $\begin{bmatrix} 0.3 & 0.7 & 0 \\ 0.1 & 0.4 & 0.5 \\ 0 & 0.2 & 0.8 \end{bmatrix}$. Is this matrix irreducible.

SQ Consider the Three States as 0, 1, 2 *

$$\begin{matrix} & 0 & 1 & 2 \\ 0 & \left[\begin{matrix} 0.3 & 0.7 & 0 \end{matrix} \right] \\ 1 & \left[\begin{matrix} 0.1 & 0.4 & 0.5 \end{matrix} \right] \\ 2 & \left[\begin{matrix} 0 & 0.2 & 0.8 \end{matrix} \right] \end{matrix}$$

In this chain we go from state 0 to state 1 with probability of 0.7 & from state 1 to state 2 with probability 0.5. Thus it is possible to go from state 0 to state 2.

\therefore The chain is irreducible & all the states are recurrent.

3. If the matrix $\begin{bmatrix} 0.4 & 0.6 & 0 & 0 \\ 0.3 & 0.7 & 0 & 0 \\ 0.2 & 0.4 & 0.1 & 0.3 \\ 0 & 0 & 0 & 1 \end{bmatrix}$ is irreducible.

$$\text{sq } P = \begin{bmatrix} 0.4 & 0.6 & 0 & 0 \\ 0.3 & 0.7 & 0 & 0 \\ 0.2 & 0.4 & 0.1 & 0.3 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$P^2 = P \cdot P = \begin{bmatrix} 0.34 & 0.66 & 0 & 0 \\ 0.33 & 0.67 & 0 & 0 \\ 0.22 & 0.44 & 0.01 & 0.33 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$P^3 = \begin{bmatrix} 0.334 & 0.666 & 0 & 0 \\ 0.333 & 0.667 & 0 & 0 \\ 0.222 & 0.444 & 0.001 & 0.33 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

W.K.T, one of the rule to prove that, a matrix is irreducible is that the sum of each row of the matrix must be equal to 1, but in matrix P^3 , the sum of third row is not equal to 1. Hence

The given matrix is not irreducible.

* If $x = [x_1, x_2, \dots, x_n]$ is the steady-state distribution of the chain whose tpm is n^{th} order square matrix 'p'.

Then $x = xp$

14) If the tpm of a markov chain is $\begin{bmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$. Find the steady state distribution of the chain.

Sol let $x = [x_1, x_2]$ is the steady-state distribution of the chain

Then $x = xp$ and $x_1 + x_2 = 1$ — ①

$$[x_1 \ x_2] = [x_1 \ x_2] \begin{bmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$[x_1 \ x_2] = \left[\frac{1}{2}x_2 \quad x_1 + \frac{1}{2}x_2 \right]$$

$$x_1 = \frac{1}{2}x_2 \Rightarrow 2x_1 = x_2 \quad \text{--- ②}$$

$$x_2 = x_1 + \frac{1}{2}x_2 \quad \text{--- ③}$$

sub eqⁿ ② in ①, we get

$$x_1 + 2x_1 = 1$$

$$3x_1 = 1$$

$$x_1 = \frac{1}{3}$$

from eqⁿ ②, we get

$$x_2 = 2 \times \frac{1}{3}$$

$$= \frac{2}{3}$$

\therefore Hence $\left[\frac{1}{3} \quad \frac{2}{3} \right]$ be the steady state distribution
of the chain.

- Q) 15) Study of passage of english text to find a vowel followed by a vowel or a consonant followed by a consonant or a vowel reveal the following transition probability matrix $\begin{bmatrix} 0.12 & 0.88 \\ 0.54 & 0.46 \end{bmatrix}$. Find the percentage of letters in the text book which are expected to be vowels.

Given that the transition probability matrix $P = \begin{bmatrix} 0.12 & 0.88 \\ 0.54 & 0.46 \end{bmatrix}$

We have to find the limiting probabilities

$$x = xp \Rightarrow x_1 + x_2 = 1 \quad \text{--- (1)}$$

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 0.12 & 0.88 \\ 0.54 & 0.46 \end{bmatrix}$$

$$x_1 = 0.12x_1 + 0.54x_2 \quad \text{--- (2)}$$

$$x_2 = 0.88x_1 + 0.46x_2 \quad \text{--- (3)}$$

Sub, $x_2 = 1 - x_1$ in (2) (from (1))

We get

$$0.12x_1 + 0.54(1 - x_1) = x_1$$

$$\Rightarrow 0.12x_1 - 0.54x_1 + 0.54 = x_1$$

$$\Rightarrow x_1 = 0.3803$$

$$\therefore x_2 = 1 - x_1$$

$$= 1 - 0.3803$$

$$x_2 = 0.6197$$

\therefore Hence, the percentage of letters in the textbook which are expected to be vowels is 38%.

- 16) The School of International studies for population found out by its survey that the mobility of the population of a state to village, town and city is the following percentage.

From \ To	Village	Town	City
Village	30%	20%	50%
Town	30%	50%	20%
City	10%	40%	50%

What will the proportion of population in village, town & city after two years? present population has proportion of 0.4, 0.3 and 0.3 village, town and city respectively. find the proportions in the long run.

Sol:- Given that the school of International studies for population found in its survey that the mobility of the population of a state to village, Town & City are.

From \ To	village	Town	city
village	30%	20%	50%
Town	30%	50%	20%
city	10%	40%	50%

which can be written as

$$P = \begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.3 & 0.5 & 0.2 \\ 0.1 & 0.4 & 0.5 \end{bmatrix}$$

The present population proportion $p_0 = [0.4, 0.3, 0.3]$

The proportion of population in village, town & city after one year is

$$\begin{aligned} p_1 &= p_0 \cdot P \\ &= [0.4 \ 0.3 \ 0.3] \begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.3 & 0.5 & 0.2 \\ 0.1 & 0.4 & 0.5 \end{bmatrix} \end{aligned}$$

$$P_1 = \begin{bmatrix} 0.24 & 0.35 & 0.4 \end{bmatrix}$$

Hence the proportion of population in village, town and city after one year

$$P_1 = \begin{bmatrix} 0.24 & 0.35 & 0.4 \end{bmatrix}$$

→ The proportion of population in village, town and city after two years

$$P_2 = P_1 \cdot P$$

$$= \begin{bmatrix} 0.24 & 0.35 & 0.4 \end{bmatrix} \begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.3 & 0.5 & 0.2 \\ 0.1 & 0.4 & 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} 0.218 & 0.387 & 0.395 \end{bmatrix}$$

Hence the proportion of population in village, town and city after two years are 0.218, 0.387 & 0.395.

(b) The proportion of population in village, town and city in long run.

i.e the steady-state distribution of the chain, then

$$x = xp \text{ and } x_1 + x_2 + x_3 = 1$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.3 & 0.5 & 0.2 \\ 0.1 & 0.4 & 0.5 \end{bmatrix}$$

$$x_1 = 0.3x_1 + 0.3x_2 + 0.1x_3$$

$$\Rightarrow 0.7x_1 = 0.3x_2 + 0.1x_3 \quad \text{--- } ①$$

$$x_2 = 0.2x_1 + 0.5x_2 + 0.4x_3$$

$$\Rightarrow 0.5x_2 = 0.2x_1 + 0.4x_3 \quad \text{--- } ②$$

$$x_3 = 0.5x_1 + 0.2x_2 + 0.5x_3$$

$$\Rightarrow 0.5x_3 = 0.5x_1 + 0.2x_2 \quad \text{--- } ③$$

$$x_1 + x_2 + x_3 = 1$$

$$x_1 = 1 - x_2 - x_3 \quad \text{--- } ④$$

Now multiplying eqⁿ ④ with 0.7 and subtracting eqⁿ ①
we get,

$$\begin{aligned} 0.7x_1 &= 0.7 - 0.7x_2 - 0.7x_3 \\ 0.7x_1 &= 0 + 0.3x_2 + 0.1x_3 \\ \hline 0 &= 0.7 - 1.0x_2 - 0.8x_3 \end{aligned}$$

$$\therefore x_2 + 0.8x_3 = 0.7 \quad \text{--- } ⑤$$

Now multiply eqⁿ ② with 0.5 and eqⁿ ③ with 0.2 then
we get:

$$\begin{aligned} 0.25x_2 &= 0.1x_1 + 0.2x_3 \\ 0.1x_3 &= 0.1x_1 + 0.04x_2 \\ \hline 0.25x_2 - 0.1x_3 &= 0.2x_3 - 0.04x_2 \\ 0.25x_2 + 0.04x_2 &= 0.2x_3 + 0.1x_3 \end{aligned}$$

$$0.29x_2 = 0.3x_3$$

$$x_2 = \left(\frac{0.3}{0.29}\right)x_3$$

$$x_2 = (1.034)x_3 \quad \text{--- (6)}$$

Now, substituting eqⁿ (6) in eqⁿ (5), we get

$$(1.034)x_3 + 0.8x_3 = 0.7$$

$$1.834x_3 = 0.7$$

$$x_3 = \frac{0.7}{1.834}$$

$$x_3 = 0.382$$

Sub $x_3 = 0.382$ in eqⁿ (6), we get

$$x_2 = (1.034)(0.382)$$

$$x_2 = 0.394$$

Put $x_2 = 0.394$ and $x_3 = 0.382$ in eqⁿ (4), we get

$$x_1 = 1 - 0.394 - 0.382$$

$$x_1 = 0.224$$

$$X = [x_1 \ x_2 \ x_3] = [0.224 \ 0.394 \ 0.382]$$

\therefore The long run proportion of population move from State to village, town and city is 22.4%, 39.4% and 38.2% respectively.

~~Ques.~~

17) If the transition probability matrix of market shares of three brands A, B and C is $\begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.35 & 0.25 & 0.4 \end{bmatrix}$ and the initial market shares are 50%, 25% and 25%. Find (i) The market shares in second and third periods. (ii) The limiting probabilities.

Sol:- Given that, the initial market shares of three brands A, B and C are 50%, 25% and 25% which can be written as $P_0 = [0.5 \ 0.25 \ 0.25]$

The transition probability matrix

$$P = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.35 & 0.25 & 0.4 \end{bmatrix}$$

The market shares in the second period = $P_0 \cdot P^2$

$$\begin{aligned} &= [0.5 \ 0.25 \ 0.25] \begin{bmatrix} 0.5050 & 0.225 & 0.27 \\ 0.435 & 0.275 & 0.29 \\ 0.48 & 0.23 & 0.29 \end{bmatrix} \\ &= [0.4813 \ 0.2387 \ 0.2800] \end{aligned}$$

Hence the market shares in the third period = $p_0 \cdot p^3$

$$= [0.5 \ 0.25 \ 0.25] \begin{bmatrix} 0.4145 & 0.2085 & 0.19 \\ 0.4675 & 0.2105 & 0.1985 \\ 0.4652 & 0.2304 & 0.221 \end{bmatrix}$$

$$= [0.4815 \ 0.2382 \ 0.2802]$$

Hence the market shares in the third period is

48.1%, 23.8%, 28.02%.

ii) The limiting probability can be interpreted as the long run market shares that can be reached to an equilibrium.

$$x = [x_1 \ x_2 \ x_3] \text{ then } x = xp \text{ and } x_1 + x_2 + x_3 = 1 \quad \text{--- (1)}$$

$$[x_1 \ x_2 \ x_3] = [x_1 \ x_2 \ x_3] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.35 & 0.25 & 0.4 \end{bmatrix}$$

$$\text{Now } x_1 = 0.4x_1 + 0.8x_2 + 0.35x_3$$

$$0.6x_1 - 0.8x_2 - 0.35x_3 = 0 \quad \text{--- (2)}$$

$$\text{Now } x_2 = 0.3x_1 + 0.1x_2 + 0.25x_3$$

$$-0.3x_1 + 0.9x_2 - 0.25x_3 = 0 \quad \text{--- (3)}$$

$$\text{Now } x_3 = 0.3x_1 + 0.1x_2 + 0.4x_3$$

Solving eqns ①, ② and ③, we get

$$x = 0.4813, \quad y = 0.2383, \quad z = 0.2804$$