

MATHEMATICAL AND STATISTICAL FOUNDATIONS

UNIT-2

Simple Linear Regression and
Correlation and Random Variables
and Probability distributions

2.0 OBJECTIVES

- To understand the importance of having simplified expression in Linear Regression
- To understand the concept of correlation
- To learn the ideas of Probability and axioms of Probability
- To understand the concept of random variables, exhaustive events.
- To classify the discrete and continuous random variables.
- To learn the method to find probability distributions.

2.1 Introduction

Linear regression is a very powerful statistical technique. Many people have some familiarity with regression just from reading the news where graphs with straight lines are overlaid on scatterplots. Linear models can be used for prediction or to evaluate whether there is linear relationship between two numerical variables. The representation is a linear

equation that combines a specific set of input values (X). As such, both the input values (X) and the output value are numeric. Linear Regression refers to a group of techniques for fitting and studying the straight line relationship between two variables. Linear regression estimates the regression coefficients β_0 and β_1 in the equation $Y_j = \beta_0 + \beta_1 X_j + \epsilon_j \in (0, 1)$ where X is the independent variable, Y is the dependent variable, β_0 is the y intercept, β_1 is the slope and ϵ is the error.

Once the intercept and slope have been estimated using least squares, various indices are studied to determine the reliability of these estimates. One of the most popular of these reliability indices is the correlation coefficient. The correlation coefficient, or simply the correlation, is an index that ranges

from -1 to 1. When the value is near zero, there is no linear relationship. As the correlation gets closer to plus or minus one, the relationship is stronger. A value of one (or negative one) indicate a perfect linear relationship between two variables. Actually, the strict interpretation of the correlation is different from that given in the last paragraph. The correlation is a parameter of the bivariate normal distribution. This distribution is used to describe the association between two variables. This association does not include a cause and effect statement. That is, the variables are not labeled as dependent and independent. One does not depend on the other. Rather, they are considered as two random variables that seem to vary together. The important point is that in linear

regression, y is assumed to be a random variable and x is assumed to be a fixed variable. In correlation analysis, both y and x are assumed to be random variable.

2.2 Simple Linear Regression

Simple linear regression, which involves two variables. One is denoted by Y , which is the variable of primary concern. Its behaviour is uncontrolled and uncertain.

That's why, Y is a random variable.

The mean of Y and the value assumed by Y depends on some other variable X . The relationship between the mean of Y and X is linear.

We know that the behaviour of Y depends on X , so Y is called the dependent variable (or) the response variable. Use observed data to develop a linear equation, express the mean of Y in terms of X . The values of X used in developing this equation is controlled. So, X is not considered as a random variable.

Select the values of x to use in our experiment and then observe the values which we assumed by the random variable y at these points.

Since the values of x don't depend on y , then x is called the "independent variable" (or) predictor variable.

Example: A physician wants to predict the concentration of a particular

drug in the blood stream (y)

based on the length of time (t)

since the drug was administered

to the patient.

Example: An economist wants to develop

an equation by which the price

of Rice in Nizamabad (y) can be

predicted from the amount of received

during the growing season in the

midwest farm belt (x).

from the above discussion γ is a random variable whose distribution depends on η . Mostly we are interested mainly in the relation between η and the mean of the corresponding distribution of the γ 's.

2.3 The Method of Least squares:

The Method of Least squares is most systematic method to fit a unique curve through the given data points and is widely used in practical computations.

Let the observed value at $x=x_i$ is y_i and corresponding value on the curve is $f(x_i)$ let e_i is the error of approximation at $x=x_i$. Then we have

$$e_i = y_i - f(x_i)$$

Consider,

$$S = \left[(y_1 - f(x_1))^2 + (y_2 - f(x_2))^2 + \dots + (y_n - f(x_n))^2 \right]$$

$$S = e_1^2 + e_2^2 + \dots + e_n^2$$

The Method of Least squares consists of minimizing S.

2.3.1 Fitting a straight line:

Let x_i, y_i (where $i=1, 2, 3, \dots, n$) be the data of an unknown function $y=f(x)$. Let $y=ax+b$ is a straight line to be fitted to the given data.

Now, minimizing S

$$\text{i.e., } \frac{\partial S}{\partial a} = 0 \text{ and } \frac{\partial S}{\partial b} = 0$$

we have,

$$S = [(y_1 - f(x_1))^2 + (y_2 - f(x_2))^2 + \dots + (y_n - f(x_n))^2]$$

$$S = [(y_1 - a - b x_1)^2 + (y_2 - a - b x_2)^2 + \dots + (y_n - a - b x_n)^2]$$

$$\text{Now, } \frac{\partial S}{\partial a} = 0$$

$$\Rightarrow 2[y_1 - a - b x_1](-1) + 2[y_2 - a - b x_2](-1) + \dots + 2[y_n - a - b x_n](-1) = 0$$

$$\Rightarrow (y_1 + y_2 + y_3 + \dots + y_n) - n a - b(x_1 + x_2 + \dots + x_n) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - n a - b \sum_{i=1}^n x_i = 0$$

$$\therefore \Rightarrow \boxed{\sum_{i=1}^n y_i = n a + b \sum_{i=1}^n x_i}$$

$$\text{Now, } \frac{\partial S}{\partial b} = 0$$

$$\Rightarrow 2[y_1 - a - b x_1](-x_1) + 2[y_2 - a - b x_2](-x_2) + \dots + 2[y_n - a - b x_n](-x_n) = 0$$

$$\Rightarrow (y_1 x_1 - a x_1 - b x_1^2) + (y_2 x_2 - a x_2 - b x_2^2) + \dots + (y_n x_n - a x_n - b x_n^2) = 0$$

$$\Rightarrow (y_1 x_1 + y_2 x_2 + \dots + y_n x_n) - a(x_1 + x_2 + \dots + x_n) - b(x_1^2 + x_2^2 + \dots + x_n^2) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i = a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2$$

\therefore The two normal equations to fit a straight line are

$$\sum y = n a + b \sum x$$

$$\sum xy = a \sum x + b \sum x^2$$

Example: Fit a straight line to the following data by the method of least squares.

x	0	1	2	3	4
y	1	1.8	3.3	4.5	6.3

Solution: Let $y = a + bx$ is the required straight line to be fitted.

The two normal equations of a straight line are

$$\sum y = na + b \sum x \quad \text{--- (1)}$$

$$\sum xy = a \sum x + b \sum x^2 \quad \text{--- (2)}$$

Here, $n = \text{no. of observations} = 5$

x	y	xy	x^2
0	1	0	0
1	1.8	1.8	1
2	3.3	6.6	4
3	4.5	13.5	9
4	6.3	25.2	16
$\sum x = 10$		$\sum y = 16.9$	$\sum xy = 47.1$
			$\sum x^2 = 30$

Substitute all the above values in equation ① & ②

$$16.9 = 5a + 10b \quad \text{--- } ③$$

$$47.1 = 10a + 30b \quad \text{--- } ④$$

By solving eqn ③ & ④,

$$\text{we get } a = 0.72$$

$$b = 1.33$$

$\therefore y = (0.72) + (1.33)x$ is the required straight line.

Example: The temperature T (in $^{\circ}\text{C}$) and length L (in mm) of a heated rod are given below, if $L = a_0 + a_1 T$. Find the best values of a_0 and a_1 ?

T	20	30	40	50	60	70
L	800.3	800.4	800.6	800.7	800.9	801.0

Solution: Let $L = a_0 + a_1 T$ is the required straight line to be fitted.

The two normal eqns of a straight line

$$\sum L = n a_0 + a_1 \sum T \quad \text{--- (1)}$$

$$\sum LT = a_0 \sum T + a_1 \sum T^2 \quad \text{--- (2)}$$

T	L	TL	T^2
20	800.3	16006	400
30	800.4	24012	900
40	800.6	32024	1600
50	800.7	40035	2500
60	800.9	48054	3600
70	801.0	56070	4900
$\sum T = 270$		$\sum L = 48039$	$\sum TL = 216201$
			$\sum T^2 = 13900$

Substitute all the above values in eqn ② & ③

$$48039 = 6a_0 + 270a_1 \quad \dots \text{---} ④$$

$$216201 = 270a_0 + 13900a_1 \quad \dots \text{---} ⑤$$

By solving eqn ④ & ⑤

we get,

$$a_0 = 799.9943$$

$$a_1 = 0.0146$$

Now we can calculate the value of a_0 & a_1 and substitute it in eqn ①

and formula of the parabola will be

$$y = 799.9943 + 0.0146x^2$$

$$(y - 799.9943) = 0.0146x^2$$

x^2	$y - 799.9943$	x	y
0.0	0.0000	0.00	0.00
1.0	0.0146	1.00	1.00
2.0	0.0584	2.00	2.00
3.0	0.1225	3.00	3.00
4.0	0.2074	4.00	4.00
5.0	0.3136	5.00	5.00

2.3.2 Fitting a Parabola (or) Second degree

Polynomial:

Let $y = ax^2 + bx + c$ - ① is a parabola to be fitted to the given data.

$$\text{WKT } S = [y_1 - a - bx_1 - cx_1^2]^2 + [y_2 - a - bx_2 - cx_2^2]^2 + \dots + [y_n - a - bx_n - cx_n^2]^2$$

To minimize S , we have $\frac{\partial S}{\partial a} = 0$,

$$\frac{\partial S}{\partial b} = 0, \quad \frac{\partial S}{\partial c} = 0.$$

$$\text{Now, } \frac{\partial S}{\partial a} = 0$$

$$2[y_1 - a - bx_1 - cx_1^2](-1) + 2[y_2 - a - bx_2 - cx_2^2](-1) + \dots + 2[y_n - a - bx_n - cx_n^2](-1) = 0$$

$$(-2)[y_1 + y_2 + \dots + y_n - na - b(x_1 + x_2 + \dots + x_n) - c(x_1^2 + x_2^2 + \dots + x_n^2)] = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - na - b \sum_{i=1}^n x_i - c \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^n y_i = na + b \sum_{i=1}^n x_i + c \sum_{i=1}^n x_i^2} - ②$$

$$\text{Now } \frac{\partial S}{\partial b} = 0$$

$$\Rightarrow 2[y_1 - a - bx_1 - cx_1^2](-x_1) + 2[y_2 - a - bx_2 - cx_2^2](-x_2)$$

$$+ \dots + 2[y_n - a - bx_n - cx_n^2](-x_n) = 0$$

$$\Rightarrow -2[y_1 x_1 - ax_1 - bx_1^2 - cx_1^3] + y_2 x_2 - ax_2 - bx_2^2 - cx_2^3 + \\ \dots + y_n x_n - ax_n - bx_n^2 - cx_n^3 = 0$$

$$\Rightarrow (y_1 x_1 + y_2 x_2 + \dots + y_n x_n) - a(x_1 + x_2 + \dots + x_n) - \\ b(x_1^2 + x_2^2 + \dots + x_n^2) - c(x_1^3 + x_2^3 + \dots + x_n^3) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 - c \sum_{i=1}^n x_i^3 = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^n x_i y_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 + c \sum_{i=1}^n x_i^3} \quad - \textcircled{3}$$

$$\text{Now } \frac{\partial S}{\partial C} = 0$$

$$\Rightarrow 2[y_1 - a - bx_1 - cx_1^2](-x_1^2) + 2[y_2 - a - bx_2 - cx_2^2](-x_2^2)$$

$$+ \dots + 2[y_n - a - bx_n - cx_n^2](-x_n^2) = 0$$

$$\Rightarrow -2[y_1 x_1^2 - ax_1^2 - bx_1^3 - cx_1^4] + y_2 x_2^2 - ax_2^2 - bx_2^3 - cx_2^4 + \\ \dots + y_n x_n^2 - ax_n^2 - bx_n^3 - cx_n^4 = 0$$

$$\Rightarrow (y_1 x_1^2 + y_2 x_2^2 + \dots + y_n x_n^2) - a(x_1^2 + x_2^2 + \dots + x_n^2) - \\ b(x_1^3 + x_2^3 + \dots + x_n^3) - c(x_1^4 + x_2^4 + \dots + x_n^4) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i x_i^0 - a \sum_{i=1}^n x_i^0 - b \sum_{i=1}^n x_i^3 - c \sum_{i=1}^n x_i^4 = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^n y_i x_i^0 = a \sum_{i=1}^n x_i^0 + b \sum_{i=1}^n x_i^3 + c \sum_{i=1}^n x_i^4} \quad \text{--- (4)}$$

\therefore The three normal equations to fit a Parabola are

$$\sum y = na + bx + cx^2$$

$$\sum xy = ax + bx^2 + cx^3$$

$$\sum x^2 y = ax^2 + bx^3 + cx^4$$

Example: Fit a Parabola of the form $y = a + bx + cx^2$ to the following data

x	1	2	3	4	5	6	7
y	2.3	5.2	9.1	16.5	29.4	35.5	54.4

By the method of least squares.

Solution: Let $y = a + bx + cx^2$ is the required Parabola to be fitted. The three normal equations to fit a Parabola

$$\sum y = na + b \sum x + c \sum x^2 \quad \text{--- (2)}$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x^3 \quad \text{--- (3)}$$

$$\sum x^2y = a \sum x^2 + b \sum x^3 + c \sum x^4 \quad \text{--- (4)}$$

x	y	xy	x^2	x^3	x^2y	x^3y	x^4
1	2.3	2.3	1	1	2.3	1	1
2	5.2	10.4	4	8	20.8	16	
3	9.7	29.1	9	27	87.3	81	
4	16.5	66	16	64	264	256	
5	29.4	147	25	125	735	625	
6	35.5	213	36	216	1278	1296	
7	54.4	380.8	49	343	2665.6	2401	
$\Sigma x = 28$		$\Sigma y = 153$	$\Sigma xy = 848.6$	$\Sigma x^2 = 140$	$\Sigma x^3y = 5053$	$\Sigma x^3 = 784$	$\Sigma x^4 = 4676$

Substitute the above values in eqn (2) (3) (4)

we get, $153 = 7a + 28b + 140c$

$$848.6 = 28a + 140b + 784c$$

$$5053 = 140a + 784b + 4676c$$

By solving the above equations

we get,

$$a = 2.3714$$

$$b = -1.0928$$

$$c = 1.1928$$

$\therefore y = 2.3714 - (1.0928)x + (1.1928)x^2$ is
the required parabola.

Example: Fit a parabola of the form
 $y = a + bx + cx^2$ to the following data

x	0	1	2	3	4
y	1	1.8	1.3	2.5	6.3

By method of least squares.

Solution: Let $y = a + bx + cx^2$ is the required parabola to be fitted. The three normal equations to fit a parabola

$$\sum y = n a + b \sum x + c \sum x^2 \quad \text{--- (1)}$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x^3 \quad \text{--- (2)}$$

$$\sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^4 \quad \text{--- (3)}$$

x	y	xy	x^2	x^3y	x^3	x^4
0	1	0	0	0	0	0
1	1.8	1.8	1	1.8	1	1
2	1.3	2.6	4	5.2	8	16
3	2.5	7.5	9	22.5	27	81
4	6.3	25.2	16	100.8	64	256
$\sum x = 10$	$\sum y = 12.9$	$\sum xy = 37.1$	$\sum x^2 = 30$	$\sum x^3y = 130.3$	$\sum x^3 = 100$	$\sum x^4 = 354$

Substitute the above values in equations ②, ③ & ④.

$$12.9 = 5a + 10b + 30c \quad \text{--- (5)}$$

$$37.1 = 10a + 30b + 100c \quad \text{--- (6)}$$

$$130.3 = 30a + 100b + 354c \quad \text{--- (7)}$$

By solving equations ⑤, ⑥ & ⑦

we get $a = 1.42$, $b = -1.07$, $c = 0.55$

$\therefore y = 1.42 - 1.07x + 0.55x^2$ is the required parabola.

2.3.3 Fitting a exponential curve

Let $y = ae^{bx}$ -① is the required curve to be fitted to the given ~~data~~ set of points i.e., (x_i, y_i)

Taking loge on both sides

$$\log_e y = \log_e(ae^{bx})$$

$$\log_e y = \log_a + \log_e e^{bx}$$

$$\log_e y = \log_a + bx \quad \text{--- ②}$$

$$\text{let } Y = \log_e y$$

$$A = \log_e a$$

$$B = b$$

$$X = x$$

Substituting all the values in eqn ②

$$Y = A + BX \text{ which is a straight line}$$

The Normal equations to fit a straight line are

$$\sum Y = NA + BX \quad \text{--- ③}$$

$$\sum XY = AX^2 + BX^3 \quad \text{--- ④}$$

By solving ③ & ④, we get A & B values.

But, $A = \log_a$

$$\Rightarrow [a = e^A]$$

$$\text{and } B = b \Rightarrow [B = b]$$

Finally, substitute a, b values in eqn ①

$$\text{i.e., } [y = ae^{bx}]$$

Example: Determine the constants a, b by the method of least squares such that $y = ae^{bx}$ fit the following table.

x	2	4	6	8	10
y	4.077	11.084	30.128	81.897	222.62

Solution: Let $y = ae^{bx}$ - ①. is the required curve to be fitted to the given set of points i.e., (x_i, y_i)

Taking loge on both sides

$$\log_e y = \log_e (ae^{bx})$$

$$\log_e y = \log_e a + \log_e^{bx}$$

$$\log y = \log a + bx - ②$$

$$\text{let } t = \log y$$

$$A = \log a$$

$$X = x$$

$$B = b$$

Substituting all the values in eqn ②

$$Y = A + BX \text{ which is a straight line.}$$

The Normal equations to fit a straight line

$$\sum Y = NA + BX - ③$$

$$\sum XY = A \sum X + BX^2 - ④$$

$x = X$	y	$y = \log y$	XY	X^2
2	4.077	1.4054	8.108	4
4	11.084	2.4055	9.6220	16
6	30.128	3.4055	20.4330	36
8	81.897	4.4055	35.244	64
10	222.62	5.4055	54.055	100
$\sum X = 30$		$\sum Y = 17.0274$	$\sum XY = 122.1648$	$\sum X^2 = 220$

Substitute all the values in eqn ③ & ④

$$17.0274 = 5A + 30B \quad \text{---} ⑤$$

$$122.1648 = 30A + 220B \quad \text{---} ⑥$$

By solving ⑤ & ⑥ we get

$$A = 0.4054, B = 0.5$$

$$\begin{aligned}\text{But } \log_e a &= A \Rightarrow a = e^A \\ &\Rightarrow a = e^{0.4054} \\ &\Rightarrow a = 1.49\end{aligned}$$

$$\begin{aligned}B &= b \Rightarrow b = B \\ &\Rightarrow b = 0.5\end{aligned}$$

Finally substitute a & b values in
eqn ①

$$\begin{aligned}\text{i.e., } y &= a e^{bx} \\ &\Rightarrow y = (1.49) e^{(0.5)x}.\end{aligned}$$

Example: Fit a curve of the form $y = a e^{bx}$

x	0	1	2	3
y	1.05	2.10	3.85	8.30

Solution: Let $y = ae^{bx}$ is the required curve
 $\quad \quad \quad -①$
 to be fitted to the given set of
 points i.e., (x_i, y_i)

Taking loge on both sides

$$\log_e y = \log_e a + \log_e e^b x \quad -②$$

$$\text{Let } Y = \log_e y$$

$$A = \log_e a$$

$$B = b$$

$$X = x$$

Substituting all the values in eqn ②

$Y = A + BX$ which is a straight line

∴ The normal equations to fit a straight line are

$$\sum Y = NA + BX \quad -③$$

$$\sum XY = A\sum X + BX^2 \quad -④$$

$x = X$	y	$Y = \log_e y$	XY	X^2
0	1.05	0.048	0	0
1	2.10	0.741	0.741	1
2	3.85	1.348	2.696	4
3	8.30	2.116	6.348	9
$\sum X = 6$		$\sum Y = 4.253$	$\sum XY = 9.785$	$\sum X^2 = 14$

Substitute all the above values in eqn.

(3) & (4)

$$4.253 = 4A + 6B \quad \text{--- (5)}$$

$$9.785 = 6A + 14B \quad \text{--- (6)}$$

By solving eqns (5) and (6)

we get,

$$A = 0.042$$

$$B = 0.6809$$

$$\text{But, } a = e^A \Rightarrow a = e^{0.042} = 1.04$$

$$\text{and } b = B \Rightarrow b = 0.6809$$

Substitute a & b values in eqn (1)

$$\text{i.e., } y = (1.04) e^{(0.6809)t}$$

2.3.4 Fitting a Power curve

Let $y = ab^x$ - (1) is the required curve
to be fitted to the given set of points
(x_i, y_i).

Taking \log_{10} on both sides

$$\log_{10}y = \log_{10}a + x \log_{10}b - (2)$$

$$\text{Let } \log_{10}y = Y, \log_{10}a = A, x = X,$$

$$\log_{10}b = B.$$

Substitute all the values in eqn (2)

$$Y = A + BX - (3)$$

which is a straight line

$$\sum Y = NA + BX - (4)$$

$$\sum XY = AX + BX^2 - (5)$$

By solving eqn (3) & (5), we get A & B values.

But $\log_{10}a = A \Rightarrow a = 10^A$ (anti log)

and $\log_{10}b = B \Rightarrow b = 10^B$ (anti log)

Substituting a & b values in eqn (1)

$$\text{i.e., } y = ab^x$$

Example:

Fit a curve of the form $y = ab^x$ to the data given below:

x	1	2	3	4	5	6
y	151	100	61	50	20	8

Solution: Let $y = ab^x \dots (1)$

Take \log_{10} on both sides

$$\log_{10} y = \log_{10}(ab^x)$$

$$\log_{10} y = \log_{10} a + \log_{10} b^x$$

$$\log_{10} y = \log_{10} a + x \log_{10} b \dots (2)$$

$$\text{Let } Y = \log_{10} y$$

$$A = \log_{10} a$$

$$X = x$$

$$B = \log_{10} b$$

Substituting all the values in (2)

$$Y = A + BX \dots (3)$$

which is a straight line

The Normal equations are

$$\sum Y = nA + BX \dots (4)$$

$$\sum XY = AX + BX^2 \dots (5)$$

$x = x$	y	$y = \log_{10} y$	xy	x^2
1	151	2.1790	2.1790	1
2	100	2	4	4
3	61	1.7853	5.3559	9
4	50	1.6990	6.7960	16
5	20	1.3010	6.5050	25
6	8	0.9031	5.4186	36
$\Sigma x = 81$		$\Sigma y = 9.8674$	$\Sigma xy = 30.2545$	$\Sigma x^2 = 91$

Here $n = 6$

Substitute all the values in eqn ④ & ⑤

$$9.8674 = 6A + 21B \quad \text{--- (6)}$$

$$30.2545 = 21A + 91B \quad \text{--- (7)}$$

By solving eqn ⑥ & ⑦

we get, $A = 2.5$, $B = -0.2427$

$$\text{But, } A = \log_{10} a \Rightarrow a = 10^A = 10^{2.5}$$

$a = 316.22$

$$\text{and } B = \log_{10} b \Rightarrow b = 10^B = 10^{-0.2427}$$

$b = 0.569$

Substitute a and b in eqn ①

$$\therefore y = (316.22)(0.569)^x$$

2.4 Properties of the Least-squares Estimators:

Under the assumptions of the simple-regression Model, the sample least-squares coefficients A and B have several properties as estimators of the population regression coefficients α and β .

→ The least-squares intercept and slope are linear estimators, that is they are linear functions of the observations y_i .

$$\text{For example, } B = \sum_{i=1}^n m_i y_i$$

$$\text{where, } m_i = \frac{x_i - \bar{x}}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

The result makes the distributions of the least-squares coefficients simple.

→ under the assumption of linearity, A and B are unbiased estimators of α and β .

$$E(A) = \alpha$$

$$E(B) = \beta$$

→ under the assumptions of linearity,
constant variance, and independence
A and B have simple sampling variances.

$$V(A) = \frac{\sigma^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$$

$$V(B) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

→ under normality, the least-squares
estimators are most efficient among
all unbiased estimators, not just among
linear estimators. This is a much
more compelling result.

→ under the assumptions, the least-squares
coefficients A and B are the maximum-
likelihood estimators of α and β .

→ under the assumption of normality, the
least-squares coefficients are themselves
normally distributed.

$$A \sim N \left[\alpha, \frac{\sigma_e^2 \sum x_i^n}{n \sum (x_i - \bar{x})^2} \right]$$

$$B \sim N \left[\beta, \frac{\sigma_e^2}{\sum (x_i - \bar{x})^2} \right]$$

Even if the errors are not normally distributed, the distributions of A and B are approximately normal, with the approximation improving as the sample size grows (The Central Limit Theorem).

2.5 : CORRELATION

Definition.

Correlation is a statistical analysis which measures and analyses the degree or extent to which two variables fluctuate with reference to each other.

The correlation expresses the relationship or independence of two sets of variables upon each other. One variable may be called the subject and the other relative (dependent).

Types of Correlation

Correlation is classified into many types.

1. Positive and negative.

2. Simple and multiple.

3. Partial and total.

4. Linear and non-linear.

1. Positive And Negative Correlation.

Positive and Negative correlation depend upon the direction of change of the variables. If the two variables tend to move together in the same direction. i.e an increase in the value of one variable is accompanied by an increase in the value of other variable.

36 or a decrease in the value of one variable is accompanied by a decrease in the value of the other variable, then the correlation is called positive or direct correlation.

Simple and multiple correlation
When we study only two variables, the relationship is described as simple correlation; example quantity of money and price level, demand and price etc. But in a multiple correlation we study more than two variables simultaneously.

Partial and Total correlation.
The study of two variables excluding some other variables is called partial correlation. for example, we study price and demand, eliminating the supply side. In total correlation, all the facts are taken into account.

Linear and Non-linear correlation.
If the ratio of change between two variables is uniform, then there will be linear correlation between them. Consider

A	2	7	12	17
B	3	9	15	21

We can see that the ratio of change between the variables is the same. If we plot these on the graph, we get a straight line.

Methods of Studying Correlation.

There are 2 different methods for finding out the relationship between variables. They are

(1) Graphic method

(1). Graphic methods

. (a) Scatter diagram or scattergram.

(b) Simple graph.

(2) Mathematical method

(a) Karl Pearson's coefficient of correlation

(b) Spearman's Rank coefficient of correlation

(c) Coefficient of concurrent deviation.

(d) Method of least squares.

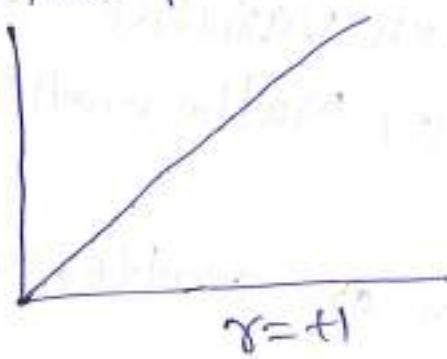
Scatter Diagram or Scattergram.

The scatter diagram is a chart obtained by plotting two variables to find out whether there is any relationship between them.

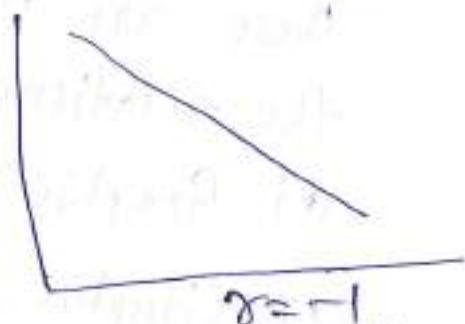
In this diagram are plotted on the horizontal axis and Y variables are plotted on the vertical axis. Thus we can know the scatter or concentration of various points.

Various scatter diagrams are briefly shown here.

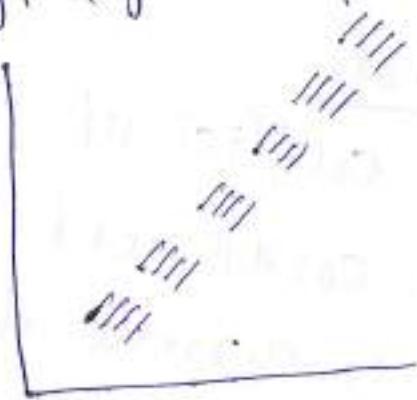
Perfect positive correlation



Perfect Negative Correlation



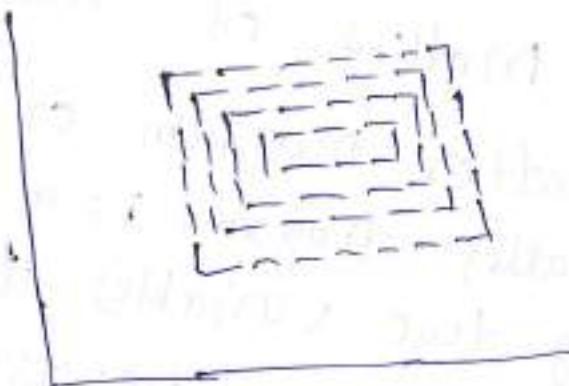
High degree of positive Correlation



High degree of negative Correlation



No correlation



2.5.1 ADVANTAGES OF SCATTER DIAGRAM

1. Scatter diagram is a simple, attractive method to find out the nature of correlation.
2. It is easy to understand.
3. A rough idea is got at a glance whether it is positive or negative correlation.

SIMPLE GRAPH

The values of the two variables are plotted on a graph paper. We get two curves, one for X variables and another for Y variables. These two curves reveal the direction and closeness of the two variables and also reveal whether are not the variables are related. If both the curves move in the same direction i.e parallel to each other, either upward or downward, correlation is said to negative.

Coefficient of correlation.

Correlation is a statistical technique used for analysing the behaviour of two or more variables. Its analysis deals with the association, between two or more variables. Statistical measures of correlation relates

w to covariation between series but not of function or causal relationship.

2.5.2 KARL PEARSON'S COEFFICIENT OF CORRELATION

Karl Pearson is a British Biometricalian and statistician suggested a mathematical method for measuring the magnitude of linear relationship between two variables. This is known as Pearsonian Coefficient of correlation. It is denoted by r . This method is most widely used. It is also called product-moment correlation coefficient.

There are several formulae to calculate r .

They are (1) $r = \frac{\text{covariance of } xy}{\sigma_x \sigma_y}$

(2) $r = \frac{\sum xy}{N \cdot \sigma_x \sigma_y}$

(3) $r = \frac{\sum xy}{\sqrt{\sum x^2} \cdot \sqrt{\sum y^2}}$

$x = (x - \bar{x})$, $y = (y - \bar{y})$ where \bar{x}, \bar{y} are mean of the series x and y .

σ_x = standard deviation of series x .

σ_y = standard deviation of series y .

Properties of Correlation Coefficient

1. Limits for correlation coefficient are $-1 \leq r \leq 1$
 Hence correlation coefficient can not exceed one numerically.

If $r=1$ correlation is perfect and positive

If $r=-1$ correlation is perfect and negative

If $r=0$ then there is no relationship between the variables.

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

If $U = \frac{X-a}{k}$ and $V = \frac{Y-b}{k}$ then $r(X, Y) = r(U, V)$

3. If X, Y are random variables and a, b, c, d are any numbers such that $a \neq 0, c \neq 0$ then

$$r(ax+b, cy+d) = \frac{ac}{|ac|} r(X, Y)$$

4. Two independent variables are uncorrelated. That is if X and Y are independent variables then $r(X, Y) = 0$.

EXAMPLES.

Example! Calculate the coefficient of correlation between age of cars and annual maintenance cost and comment:

42

Age of cars

2 4 6 7 8 10 12

Annual maintenance cost 1600 1500 1800 1900 1700 2100 2000

Solution: Let age of cars = x .Annual maintenance = y .

Computation of coefficient of correlation

$$x = x - \bar{x} \quad y = (y - \bar{y}) / 100 \quad y^2 \quad \Sigma x^2$$

$$\bar{x} = 7 \quad \bar{y} = 1800$$

x	y	$x - \bar{x}$	$y - \bar{y}$	x^2	y^2	xy	Σx^2	Σy^2
2	1600	-5	-2	25	4	10		
4	1500	-3	3	9	9	9		
6	1800	-1	0	1	1	0		
7	1900	0	1	0	1	1		
8	1700	1	-1	1	9	9		
10	2100	3	3	9	81	27		
12	2000	5	2	25	4	10		

$\Sigma x = 49$

$\Sigma y = 12600$

$\Sigma x^2 = 90$

$\Sigma y^2 = 28$

$$\gamma = \frac{\Sigma xy}{\sqrt{\Sigma x^2 \Sigma y^2}} = \frac{37}{\sqrt{90 \times 28}} = 0.836$$

We can observe that there is a high degree of positive correlation between age of cars and annual maintenance cost.

Example: find Karl Pearson's coefficient of Correlation from the following data

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

Solution: Computation of Karl Pearson's coefficient of correlation

Wages (x)	$x = x - \bar{x}$	x^2	Cost of living (y)	$y = y - \bar{y}$	y^2	X
	$\bar{x} = 99$			$\bar{y} = 95$		
100	+1	1	98	+3	9	+
101	+2	4	99	+4	16	+
102	+3	9	99	+4	16	+
102	+3	9	97	+2	4	+
100	+1	1	95	0	0	-
99	0	0	92	-3	9	-
97	-2	4	95	0	0	0
98	-1	1	94	-1	1	-
96	-3	9	90	-5	25	-
95	-4	16	91	-4	16	+
	$\sum x = 0$	$\sum x^2 = 54$	$\sum y = 950$	$\sum y^2 = 5184$	$\sum xy = 0$	$\Sigma y = 96$
$\Sigma n = 990$						

$$\therefore r = \frac{\sum xy}{\sqrt{\sum x^2 \times \sum y^2}} = \frac{61}{\sqrt{54 \times 96}} = \frac{61}{\sqrt{5184}} = \frac{61}{72} = 0.847$$

$$\therefore r = +0.847$$

w Example! Psychological tests of intelligence and of engineering ability were applied to 10 students. Here is a record of ungrouped data showing intelligence ratio (I.R) and Engineering ratio (E.R) calculate the coefficient of correlation.

student	A	B	C	D	E	F	G	H	I	J
I.R	105	104	102	101	100	99	98	96	93	92
E.R	101	103	100	98	95	96	104	92	97	94

Solution! We construct the following table

Student	Intelligence ratio		Engineering ratio		x^2	y^2	xy
	n	$\bar{x} = X$	y	$\bar{y} = Y$			
A	105	6	101	3	36	9	18
B	104	5	103	5	25	25	2
C	102	3	100	2	9	4	6
D	101	2	98	0	4	0	0
E	100	1	95	-2	1	4	-2
F	99	0	96	-2	0	36	-18
G	98	-1	104	6	1	36	36
H	96	-3	92	-6	9	36	-36
I	93	-6	97	-1	36	1	16
J	92	-7	94	-4	49	16	16
Totals	990	0	980	0	170	140	

$$= \frac{322}{\sqrt{196 \times 588}} = \frac{322}{359.48} = 0.95.$$

45

Example: find if there is any correlation between the heights and weight given below.

Height in inches	57	59	62	63	64	65	55	58	51
Weight in lbs	113	117	126	126	130	129	111	116	11

Solution: Coefficient of correlation $\gamma = \frac{\sum XY}{\sqrt{\sum X^2 \times \sum Y^2}}$

Height in inches x	Computation of coefficient of correlation					
	Deviation from Mean (60) $x = x - \bar{x}$	Square of deviation x^2	Weight y	Deviation square of mean deviation $y = y - \bar{y}$	y^2	
57	-3	9	113	7	49	
59	-1	1	117	-3	9	
62	2	4	126	6	36	
63	3	9	126	6	36	
64	4	16	130	10	100	
65	5	25	129	9	81	
55	-5	25	111	-9	81	
58	-2	4	116	-4	16	
57	-3	9	112	-8	64	
540	0	102	1080	0	0	

$$\gamma = \frac{216}{\sqrt{102 \times 471}} = 0.98$$

Example: Calculate coefficient of Correlation from the following data

X	12	9	8	10	11	13	7
Y	14	8	6	9	11	12	3

Solution: In both series items are in small number. So there is no need to take deviations.

We use the formula $r = \frac{\text{Covariance of } XY}{\sigma(X) \cdot \sigma(Y)}$

Computation of coefficient of correlation:

X	Y	X^2	Y^2	XY
12	14	144	196	168
9	8	81	64	72
8	6	64	36	48
10	9	100	81	90
11	11	121	121	121
13	12	169	144	156
7	3	49	9	21
$\Sigma X = 70$		$\Sigma Y = 63$	$\Sigma X^2 = 728$	$\Sigma Y^2 = 651$
				$\Sigma XY = 646$

$$r = \frac{(\Sigma XY \times N) - (\Sigma X \times \Sigma Y)}{\sqrt{(\Sigma X^2 \times N - (\Sigma X)^2) \cdot (\Sigma Y^2 \times N - (\Sigma Y)^2)}}$$

Here $N = 7$

$$r = \frac{(646 \times 7) - (70 \times 63)}{\sqrt{728 \times 7 - 70^2} \cdot \sqrt{651 \times 7 - 63^2}} = \frac{4732 - 4410}{\sqrt{5096 - 4900} \cdot \sqrt{4555}} = \frac{322}{\sqrt{196} \cdot \sqrt{39}} = \frac{322}{14 \cdot 6.24} = \frac{322}{87.36} = 0.37$$

from this table, mean of x , i.e. $\bar{x} = \frac{990}{10} = 99$,

mean of y i.e. $\bar{y} = \frac{980}{10} = 98$

$\sum x^2 = 170$; $\sum y^2 = 140$ and $\sum xy = 92$

substituting these values in $r = \frac{\sum xy}{\sqrt{(\sum x^2)(\sum y^2)}} =$

$$\frac{92}{\sqrt{170 \times 140}} = \frac{92}{154.3} = 0.59.$$

2.6 WHEN DEVIATIONS ARE TAKEN FROM AN ASSUMED MEAN.

When actual mean is not a whole number, but a fraction or when the series is large, the calculation by direct method will involve a lot of time. To avoid such tedious calculation, we can use the assumed mean method.

$$\sum xy - \frac{\sum x \sum y}{n}$$

∴ formula =

$$\frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{\left(\sum x^2 - \frac{(\sum x)^2}{n} \right) \sum y^2 - \frac{(\sum y)^2}{n}}}$$

where

x = deviation of the items of n -series from assumed mean. i.e. $x = x - A$

y = deviation of the items of y -series from an assumed mean. i.e $y = (y - A)$.

N = Number of items

Σxy = The total of the product of the deviations of x and y -series from their assumed mean.

Σx^2 = The total of the squares of the deviations of x -series from an assumed mean.

Σy^2 = The total of the squares of the deviation of y -series from an assumed mean.

$\Sigma x \cdot \Sigma y$ = Total of the deviation of x -series for assumed mean.

Σy = Total of the deviation of y -series for assumed mean.

Example: Find the coefficient of correlation between x and y

x	1	2	3	4	5	6	7	8	9
y	12	11	13	15	14	17	16	19	18

Solution:

X	Y	X^2	Y^2	XY
1	12	1	144	12
2	11	4	121	22
3	13	9	169	39
4	15	16	225	60
5	14	25	196	70
6	17	36	289	102
7	16	49	256	112
8	19	64	361	132
9	18	81	324	162
45	135	285	2085	731

$$\bar{X} = \frac{45}{9} = 5 ; \quad \bar{Y} = \frac{135}{9} = 15$$

∴ Using the formula $r = \frac{\text{Cov}(X,Y)}{\sigma_X \sigma_Y}$

$$= \frac{\sum XY - n(\bar{X})(\bar{Y})}{\sqrt{(\sum X^2 - n(\bar{X})^2)(\sum Y^2 - n(\bar{Y})^2)}}$$

$$= \frac{731 - 9(5)(15)}{\sqrt{(285 - 9 \times 25)(2085 - 9 \times 225)}}$$

$$= 0.933$$

Example: Calculate Karl Pearson's Correlation coefficient for the following paired data.

X 28 41 40 38 35 33 40 32 36 3
Y 23 34 33 34 30 26 28 31 36 3

what inference would you draw from the result?

Solution: Computation of correlation coefficient

n	Deviation from assumed mean $x = n - 35$ $(x = x - \bar{x})$	Square of deviation x^2	y	Deviation from assumed mean $y = y - 31$ $(y = y - \bar{y})$	Square of deviation y^2	Product xy
28	-7	49	23	-8	64	18
41	+6	36	34	+3	9	10
40	+5	25	33	+2	4	9
38	+3	9	34	+3	9	0
35	0	0	30	-1	1	10
33	-2	4	26	-5	25	-1
40	+5	25	28	-3	9	0
32	-3	9	31	0	25	5
35	+1	1	36	+5	0	49
33	-2	4	38	+7	49	8
$\Sigma x = 355$		$\Sigma x^2 = 6$		$\Sigma y = 313$	$\Sigma y^2 = 195$	

$$N=10, \text{ take } n=35 \text{ and } \bar{y}=31$$

Applying to the above data, the formula

$$\therefore r = \frac{\frac{\sum xy - \frac{\sum x\bar{y}}{N}}{\sqrt{\left[\sum x^2 - \frac{(\sum x)^2}{N} \right] \times \left[\sum y^2 - \frac{(\sum y)^2}{N} \right]}}}{= \frac{\frac{49 - \frac{6 \times 3}{10}}{\sqrt{162 - \frac{6^2}{10}} \cdot \sqrt{195 - \frac{7^2}{10}}}}{\sqrt{1584} \sqrt{194.9}} = \frac{77.2}{= 0.45}}$$

Example: Calculate the correlation coefficient for the following heights of fathers (X) and their sons (Y):

X	65	66	67	67	68	69	70	72
---	----	----	----	----	----	----	----	----

Y	67	68	65	68	72	72	69	71
---	----	----	----	----	----	----	----	----

Solution: Computation of correlation coefficient

X	Y	x^2	y^2	xy
65	67	4225	4489	4355
66	68	4356	4624	4488
67	65	4489	4225	4355
67	68	4489	5184	4556
68	72	4624	5184	4896
69	72	4761	5184	4986
70	69	4900	5041	4830
72	71	5184	5184	5112
Sum		$\Sigma x^2 = 37028$	$\Sigma y^2 = 38132$	$\Sigma xy = 37560$

$$\therefore \bar{x} = \frac{1}{n} \sum x = \frac{544}{8} = 68, \bar{y} = \frac{1}{n} \sum y = \frac{552}{8} = 69$$

$$\gamma = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\frac{1}{n} \sum xy - (\bar{x})(\bar{y})}{\sqrt{\left(\frac{1}{n} \sum x^2 - (\bar{x})^2\right) \left(\frac{1}{n} \sum y^2 - (\bar{y})^2\right)}} \\ = \frac{\frac{1}{8} \sum xy - (68)(69)}{\sqrt{\left(\frac{37028}{8} - (68)^2\right) \left(\frac{38132}{8} - (69)^2\right)}} = 0.603.$$

Example: find a suitable coefficient of Correlation for the following data:

fertiliser used	15	18	20	24	30	35	40	50
productivity	85	93	95	105	120	130	150	1

Solution: Computation of coefficient of correlation

	Fertiliser used (X)	X̄	X²	Y	Ȳ	Y - Ȳ	Y²	XY
x	$x = x - \bar{x}$ $= x - 29$					$y = y - \bar{y}$ $= y - 119$		
15	-14	196	85			-34	1156	496
18	-11	121	93			-26	696	286
20	-9	81	95			-24	576	216
24	-5	25	105			-14	196	90
30	1	1	120			1	1	1
35	6	36	130			11	121	66
40	11	121	150			31	961	341
50	21	441	160			41	1681	861
		$\sum x = 14$	$\sum x^2 = 1000$			$\sum y = 14$	$\sum y^2 = 5368$	$\sum xy = 21$

$$\therefore N=8, \quad r = \frac{\sum xy - (\bar{x})(\bar{y})}{\sqrt{(\bar{x})^2 N - (\bar{x})^2} \sqrt{(\bar{y})^2 N - (\bar{y})^2}}$$

Example: find out the coefficient of correlation in the following case.

following cast.	65	66	67	67	68	69	71	72
height of father								
height of son	67.	68	64	68	72	70	69	73

A coefficient of Correlation 01

Solution: Computation of coefficient of Correlation

x	$x = n - \bar{x}$	x^2	y	$y = y - \bar{y}$	y^2	xy
65	-2	4	67	-1	1	2
66	-1	1	68	0	0	0
67	0	0	64	-4	16	0
67	0	0	68	0	0	0
67	0	0	72	4	16	0
68	1	1	70	2	4	0
69	2	4	69	1	1	0
71	4	16	70	2	4	0
73	6	36				0
$\Sigma x = 516$	$\Sigma x = 10$	$\Sigma x^2 = 62$	$\Sigma y = 548$	$\Sigma y = 4$	$\Sigma y^2 = 42$	

$$\Sigma n = 546 \quad \Sigma k = 10$$

$$\sum k=10$$

$$\Sigma k^2 = 62$$

$$S_{\text{eff}} = 548$$

$$\Sigma Y = 4 \quad \Sigma Y^2 = 42$$

$$N = \text{no. of items} = 8$$

$$\begin{aligned}
 \text{Coefficient of Correlation } r &= \frac{\sum xy - \frac{(\sum x)(\sum y)}{N}}{\sqrt{\left[\sum x^2 - \frac{(\sum x)^2}{N} \right] \left[\sum y^2 - \frac{(\sum y)^2}{N} \right]}} \\
 &= \frac{(\sum xy)N - (\sum x)(\sum y)}{\sqrt{[(\sum x)^2 N - (\sum x)^2] \times [(\sum y)^2 N - (\sum y)^2]}} \\
 &= \frac{(26 \times 8) - (10 \times 4)}{\sqrt{(62 \times 8 - 10^2) \times (42 \times 8 - 4^2)}} \\
 &= \frac{168}{355.98} \\
 &= 0.472
 \end{aligned}$$

2.7 RANK CORRELATION COEFFICIENT

A British Psychologist Charles Edward Spearman found out the method of finding the coefficient of correlation by ranks. This method is based on rank and is useful in dealing with qualitative such as morality, character, intelligence and beauty. It can not be measured quantitatively as in the case of Pearson's coefficient of correlation. It is based on the ranks given to the observations. Rank correlation is applicable only to the individual observations.

$$\rho = 1 - \frac{6 \sum D^2}{N(N^2-1)}$$

where $\rho \rightarrow$ Rank coefficient of Correlation
 $D^2 \rightarrow$ Sum of the squares of the differences
 of two ranks.

$N \rightarrow$ Number of paired observations.

2.7.1 PROPERTIES OF RANK CORRELATION COEFFICIENT

1. The value of ρ lies between +1 and -1.
2. If $\rho=1$, there is complete agreement in the order of the ranks and the direction of the rank is same.
3. If $\rho=-1$ then there is complete disagreement in the order of the ranks and they are in opposite directions.

Procedure to solve problems:

1. When the ranks are given.
 Step-1. Compute the difference of two ranks and denote it by D .

Step-2. Square D and get $\sum D^2$.

Step-3. obtain ρ by substituting the figures in the formula.

2. When the ranks are not given, but actual

So if data are given, then we must give ranks. We can give ranks by taking the highest as 1 or the lowest value as 1, next to the highest as 2 and follow the same procedure for both the variables.

Example: Ten competitors in a musical test were ranked by the three judges A, B and C in the following order.

Ranks by A 1 6 5 10 3 2 4 9 7 8

Ranks by B 3 5 8 4 7 10 2 1 6 9

Ranks by C 6 4 9 8 1 2 3 10 5 7

Using rank correlation method, discuss which pair of judges has the nearest approach to common likings in music.

Solution: Here $N=10$.

	X	Y	Z	$D_1 = X - Y$	$D_2 = Y - Z$	$D_3 = Z - X$	D_1^2	D_2^2	D_3^2
	1	3	6	-2	-5	-3	4	25	9
	6	5	4	1	2	1	1	4	1
	5	8	9	-3	-4	-1	9	16	1
	10	4	8	6	2	-4	36	4	1
	3	7	1	-4	2	6	16	4	0
	2	10	2	-8	0	8	64	4	1
	4	2	3	2	1	-9	64	1	1
	9	1	10	8	-1	1	1	1	4
	7	6	5	1	2	2	1	1	1
	8	9	7	-1	-1	-2	1	1	1
Total				$\sum D_1 = 0$	$\sum D_2 = 0$	$\sum D_3 = 0$	$\sum D_1^2 = 200$	$\sum D_2^2 = 60$	$\sum D_3^2 = 60$

$$\rho_1(X, Y) = 1 - \frac{6 \sum D_1^2}{N(N^2-1)} = 1 - \frac{6 \times 200}{10 \times 99} = 1 - \frac{40}{33} = \frac{-7}{33}$$

$$\rho_2(X, Z) = 1 - \frac{6 \sum D_2^2}{N(N^2-1)} = 1 - \frac{6 \times 60}{10 \times 99} = 1 - \frac{4}{11} = \frac{7}{11}$$

$$\rho_3(Y, Z) = 1 - \frac{6 \sum D_3^2}{N(N^2-1)} = 1 - \frac{214}{165} = -\frac{49}{165}$$

Since $\rho_2(X, Z)$ is maximum, we conclude that the pair of judges A and C has the nearest approach to common likings in music.

Example! The ranks of 60 students in mathematics and statistics are as follows. (1,1), (2,10), (3,3), (4,4), (5,5), (6,7), (7,2), (8,6), (9,8), (10,11), (11,15), (12,9), (13,14), (14,12), (15,16), (16,13). Calculate the rank correlation coefficient for proficiencies of this group in mathematics and statistics.

Solution:

Ranks in Maths (X)

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

Ranks in Statistics (Y)

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

$$D = X - Y$$

$$D = 0 - 8 + 1 + 2 - 13 - 1 + 5 - 2 + 1 - 1 = 0$$

$$D^2$$

$$D^2 = 0^2 + 64 + 0^2 + 0^2 + 1^2 + 25 + 4 + 1 + 1 = 136$$

∴ Rank correlation coefficient

$$r = 1 - \frac{6 \sum d^2}{N(N^2-1)} = 1 - \frac{6 \times 136}{16 \times 255} = 1 - \frac{1}{5}$$

$$= \frac{4}{5} = 0.8.$$

Example: A random sample of 5 college students is selected and their grades in Mathematics and Statistics are found to be

	1	2	3	4	5
Mathematics	85	60	73	40	90
Statistics	93	75	65	50	80

Calculate Pearman's rank correlation coefficient

Solution:

Marks in Mathematics X	Ranks x	Marks in Statistics Y	Rank y	Rank difference $n-y$	D^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

$\sum D^2 = 4$

Here $N=5$, $\sum d^2=4$

* Pearman's Rank correlation

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

$$= 1 - \frac{6 \times 4}{5(5^2-1)} = 1 - \frac{24}{5 \times 24}$$

$$= 1 - \frac{1}{5} = 1 - 0.2 = 0.8.$$

Example: following are the rank obtained by 10 students in two subjects, statistics and mathematics. To what extent the knowledge of the students in two subjects is related?

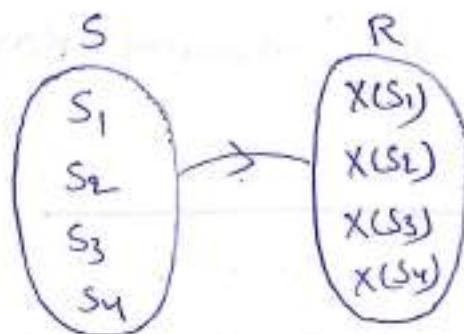
	1	2	3	4	5	6	7	8	9	10
Statistics	1	2	3	4	5	6	7	8	9	10
Mathematics	2	4	1	5	3	9	7	10	6	8

(x)	Rank in Statistics	Rank in mathy (y)	$D = (x-y)$	D^2	
				1	4
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	

$$\begin{aligned}
 l &= 1 - \frac{6\varepsilon D^2}{N(n^2-1)} \\
 &= 1 - \frac{6 \times 40}{10(10^2-1)} \\
 &= 1 - \frac{240}{10(100-1)} \\
 &= 1 - \frac{240}{990} \\
 &= 1 - 0.24 \\
 &= 0.76.
 \end{aligned}$$

2.8 Random Variable:

Definition: A random variable X on a sample space S is a function $X: S \rightarrow R$ that assigns a real number ($X(s)$) to each sample point $s \in S$.



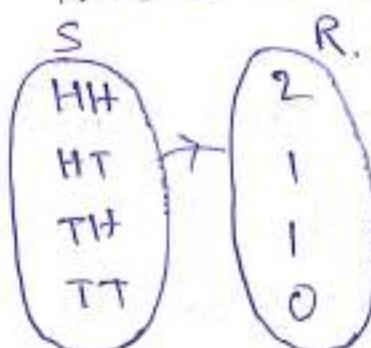
Note: Set of all possible values of a random variable X is called Range of X .

Example: In tossing two coins

Sample space $S = \{HH, HT, TH, TT\}$

Define a random variable $X: S \rightarrow R$ by

$$X(s) = \text{"No. of Heads"}$$



Range of $X = \{2, 1, 0\}$ a set of real numbers.

$\therefore X$ is a random variable

There are two types of Random Variables.

- 1) Discrete Random Variables
- 2) Continuous Random Variables

1) Discrete Random Variables:-

A random variable X which can take only a finite number of discrete values in an interval of domain is called a discrete random variable.

(or)

If the random variable takes the values only on the set $\{0, 1, 2, \dots, n\}$ is called a discrete random variable.

Example: Tossing a coin, throwing a die,

The no. of defective items in a sample of electric bulbs, the no. of printing mistakes in each page of a book,

The no. of telephone calls received by the telephone operator, no. of Accidents in a street etc.

2) Continuous Random Variable:

A random variable X which can take values continuously i.e., which takes all possible values in a given interval is called a continuous random variable.

Example: The height, age and weight of individuals, temperature and time, life span of an electric bulb.

Depending on random variable we obtain the following probability distributions.

Probability distribution function.

Let X be a random variable. Then the probability distribution function associated with X is defined as the probability that the outcome of an

experiment will be one of the outcomes
for which $X(s) \leq x, \forall s \in S$

$$\text{i.e., } F_X(x) = P(X \leq x) = P\{s : X(s) \leq x\}, \\ -\infty < x < \infty$$

is called the distribution function of X .

2.8.1 Discrete Probability distribution (probability Mass function):

Let X be a discrete random variable.
Then the probability distribution given by
 X is known as discrete probability distribution
which is represented by the following
probability mass function $p_i = P(X=x_i) = P(x_i)$
for $i=1, 2, 3, \dots$

with the properties (i) $P(x_i) \geq 0, \forall i$

$$(ii) \sum_{i=1}^{\infty} p_i = 1, i=1, 2, 3, \dots$$

2.8.2 Cumulative Distribution Function of a Discrete Random Variable:

Let X is a discrete Random Variable.

Then the Discrete distribution function
 (or) cumulative Distribution function
 $F(x)$ is defined by $F(x) = P(X \leq x)$
 $= \sum_{i=1}^x P(X_i)$

where, x is any integer.

Example: In tossing two coins,

Then $S = \{HH, HT, TH, TT\}$

Define $X: S \rightarrow R$ by $X(S) = \text{no. of heads}$.

Range of $X = \{0, 1, 2\}$

$X=x$	0	1	2
$P(X=x)$	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{4}$

is called P.m.f
 (D.P.D)

$X=x$	0	1	2
$P(X \leq x)$	$\frac{1}{4}$	$\frac{3}{4}$	1

is called Distribution
 function (or) cumulative
 distribution function.

2.8.3 Probability Density function:

The Probability density function $f_X(x)$ is defined as the derivative of the Probability distribution function, $F_X(x)$ of the random variable X .

$$\text{Thus, } f_X(x) = \frac{d}{dx}[F_X(x)]$$

Properties:

i) If F is the distribution function of a random variable X and if $a < b$, then

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

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

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

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

Note:

If $P(X=a) = P(X=b) = 0$ then

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

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

2) All distribution functions are monotonically increasing and lie b/w '0' and '1'
 i.e., if F is the distribution function of the random variable X , then

$$(i) 0 \leq F(x) \leq 1$$

$$(ii) F(x) < F(y) \text{ when } x < y$$

3) (i) $F(-\infty) = \lim_{x \rightarrow -\infty} F(x) = 0$

(ii) $F(\infty) = \lim_{x \rightarrow \infty} F(x) = 1.$

2.8.4 Continuous Probability distribution:

Let X be a continuous random variable.
 Then the probability distribution given by X is known as continuous probability distribution which is represented by $f(x)$ and known as probability density function (P.d.f) with the properties

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

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

2.8.5 Cumulative Distribution Function of a continuous random Variable:

The cumulative distribution (or) simply the distribution function of a continuous random variable X is denoted by $F(x)$ and is defined as

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(x) dx$$

Properties:

- 1) $0 \leq F(x) \leq 1$, $-\infty < x < \infty$
- 2) $F'(x) = f(x) \geq 0$, so that $F(x)$ is a non-decreasing function
- 3) $F(-\infty) = 0$
- 4) $F(\infty) = 1$
- 5) $F(x)$ is a continuous function of x on the right
- 6) The discontinuities of $F(x)$ are countable
- 7) $P(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$
- 8) since, $F'(x) = f(x)$, we have $\frac{d}{dx}[F(x)] = f(x)$
 $\Rightarrow dF = f(x) dx$

This is known as probability differential of X .

2.9 General Properties of a Discrete Probability Distribution:

Mean:

The mean of a discrete probability distribution is the mathematical expectation of a discrete random variable.

If X is a discrete random variable which takes the values x_1, x_2, \dots, x_K and whose probabilities are $f(x_1), f(x_2), \dots, f(x_K)$, then its mean (or) mathematical expectation (or) expected value is $\sum_{i=1}^K x_i f(x_i)$. The mean is

denoted by μ .

$$\therefore \mu = \sum_{i=1}^K x_i f(x_i)$$

Example: Let X be a discrete random variable taking the values 1, 2, ..., 6 with probability $f(x_i) = \frac{1}{6}$

$$\text{Then } \mu = \sum_{i=1}^6 x_i f(x_i)$$

$$= \frac{1}{6}(1+2+\dots+6)$$

$$\mu = \frac{7}{2}$$

Variance:

Let X be a discrete random variable.
Then the variance of X , denoted by σ^2 is,

$$\sigma^2 = \sum (x - \mu)^2 f(x)$$

(or)

$$\sigma^2 = \sum x^2 f(x) - \mu^2$$

Example: Determine the variance of the probability distribution of the number of points rolled with a die.

Solution: since $f(x) = \frac{1}{6}$

$$\text{we get } \mu = \frac{7}{2}$$

$$\text{and } \sum x^2 f(x) = (1^2 + 2^2 + 3^2 + 4^2 + 5^2 + 6^2) \frac{1}{6} = \frac{91}{6}$$

$$\text{Hence } \sigma^2 = \frac{91}{6} - \left(\frac{7}{2}\right)^2 = \frac{35}{12}$$

Standard Deviation:

Let X be a discrete random variable. The standard deviation of a discrete random variable X is denoted by σ , is the positive square root of σ^2

$$\text{i.e., } \sigma = \sqrt{\sum (x - \mu)^2 f(x)}$$

$$\sigma = \sqrt{\sum x^2 f(x) - \mu^2}$$

Example: If the variance $\sigma^2 = 2.0979$

$$\begin{aligned} \text{Then standard deviation } \sigma &= \sqrt{2.0979} \\ &= 1.448 \end{aligned}$$

Results:

- 1) Variance of constant is zero
- 2) $V(KX) = K^2 V(X)$, where K -constant
- 3) If X is a random variable and K is a constant, then $V(X+K) = V(X)$
- 4) If X is a discrete random variable, then $V(ax+b) = a^2 V(X)$, where $V(X)$ is a variance of X and a, b are constants.

5) If X and Y are two independent random variables, then $V(X+Y) = V(X) + V(Y)$

Example: A random variable X has the following probability function:

X	0	1	2	3	4	5	6	7
$P(X)$	0	K	$2K$	$2K$	$3K^2$	K^2	$2K^2$	$7K^2 + K$

- (i) Find K
- (ii) Evaluate $P(X \leq 6)$, $P(X \geq 6)$
- (iii) Evaluate $P(0 < X \leq 5)$, $P(0 \leq X \leq 4)$
- (iv) Find Mean
- (v) Find Variance

Solution: (i) we know that $\sum_{x=0} P(X=x) = 1$

$$\text{i.e., } 0 + K + 2K + 2K + 3K^2 + K^2 + 2K^2 + 7K^2 + K = 1$$

$$\text{i.e., } 10K^2 + 9K - 1 = 0$$

$$\text{i.e., } (10K - 1)(K + 1) = 0$$

$$\therefore K = \frac{1}{10} = 0.1 \quad (\because K \neq -1)$$

$$\begin{aligned}
 \text{(ii) } P(X < 6) &= P(X=0) + P(X=1) + \dots + P(X=5) \\
 &= 0 + K + 2K + 2K + 3K + K^2 \\
 &= 8K + K^2 \\
 &= 0.8 + 0.01 \\
 &= 0.81
 \end{aligned}$$

$$\begin{aligned}
 P(X \geq 6) &= 1 - P(X < 6) \\
 &= 1 - 0.81 \\
 &= 0.19
 \end{aligned}$$

$$\begin{aligned}
 \text{(iii) } P(0 < X < 5) &= P(X=1) + P(X=2) + P(X=3) + P(X=4) \\
 &= K + 2K + 2K + 3K \\
 &= 8K = \frac{8}{10} = 0.8 \\
 P(0 \leq X \leq 4) &= P(X=0) + P(X=1) + P(X=2) + \\
 &\quad P(X=3) + P(X=4) \\
 &= 0 + K + 2K + 2K + 3K \\
 &= 8K = 8\left(\frac{1}{10}\right) = 0.8
 \end{aligned}$$

$$\begin{aligned}
 \text{(iv) Mean } (\bar{x}) &= \sum_{i=0}^7 f_i x_i \\
 &= 0(0) + 1(1) + 2(2K) + 3(2K) + 4(3K) \\
 &\quad + 5(K^2) + 6(2K^2) + 7(7K^2 + K)
 \end{aligned}$$

$$= 66K^2 + 30K$$

$$= \frac{66}{100} + \frac{30}{10} = 0.66 + 3 = 3.66 \quad [\because K = \frac{1}{10}]$$

v) Variance = $\sum_{i=0}^{7} P_i f_i^2 - \bar{x}^2$

$$= K + 8K + 18K + 48K + 25K^2 + 72K^2 + \\ 343K^2 + 49K - (3.66)^2$$

$$= 440K^2 + 124K - (3.66)^2$$

$$= \frac{440}{100} + \frac{124}{10} - (3.66)^2$$

$$= 4.4 + 12.4 - 13.3956$$

$$= 3.4044.$$

Example: Two dice are thrown. Let X assign to each point (a_1, b) ~~is~~ S the maximum of its numbers i.e., $X(a_1, b) = \max(a_1, b)$. Find the probability distribution. X is a random variable with $X(S) = \{1, 2, 3, 4, 5, 6\}$. Also find the mean and variance of the distribution.

Solution: The total no. of cases are = 36

The maximum number could be 1, 2, 3, 4, 5, 6
 i.e., $X(S) = X(a_1 b) = \max(a_1, b)$.

$$P(1) = P(X=1) = P(1, 1) = \frac{1}{36}$$

$$P(2) = P(X=2) = \frac{3}{36} \quad (\because \text{for maximum } 2)$$

$$P(3) = P(X=3) = \frac{5}{36} \quad (\because \text{for max } 3)$$

$$P(4) = P(X=4) = \frac{7}{36} \quad (\because \text{for max } 4)$$

$$P(5) = P(X=5) = \frac{9}{36} \quad (\because \text{for max } 5)$$

$$P(6) = P(X=6) = \frac{11}{36} \quad (\because \text{for max } 6)$$

\therefore The required probability distribution is

X	1	2	3	4	5	6
$P(X)$	$\frac{1}{36}$	$\frac{3}{36}$	$\frac{5}{36}$	$\frac{7}{36}$	$\frac{9}{36}$	$\frac{11}{36}$

$$\begin{aligned}
 \text{(i) Mean } (\bar{x}) &= \sum_{i=1}^6 p_i x_i \\
 &= 1 \cdot \frac{1}{36} + 2 \cdot \frac{3}{36} + 3 \cdot \frac{5}{36} + 4 \cdot \frac{7}{36} + 5 \cdot \frac{9}{36} \\
 &\quad + 6 \cdot \frac{11}{36}
 \end{aligned}$$

$$= \frac{1}{36} (1+6+15+28+45+66)$$

$$= \frac{161}{36} = 4.47$$

(ii) Variance (σ^2) = $\sum_{i=1}^6 p_i x_i^2 - \bar{x}^2$

$$= \frac{1}{36}(1)^2 + \frac{3}{36}(2)^2 + \frac{5}{36}(3)^2 + \frac{7}{36}(4)^2 + \frac{9}{36}(5)^2 + \frac{11}{36}(6)^2 - (4.47)^2$$

$$= \frac{1}{36}(1+12+45+112+225+396) - (4.47)^2$$

$$= \frac{791}{36} - 19.9808$$

$$= 21.97 - 19.9808 = 1.9912.$$

Example: Find the mean and variance of the uniform probability distribution given

$$\text{by } f(x) = \frac{1}{n} \text{ for } x = 1, 2, 3, \dots, n.$$

Solution: The probability distribution is

x	1	2	3	n
$f(x)$	$\frac{1}{n}$	$\frac{1}{n}$	$\frac{1}{n}$	$\frac{1}{n}$

$$\begin{aligned}
 \text{(i) Mean } (\bar{x}) &= \sum_{i=1}^n x_i f(x_i) \\
 &= 1 \cdot \frac{1}{n} + 2 \cdot \frac{1}{n} + \dots + n \cdot \frac{1}{n} \\
 &= \frac{1}{n} (1+2+3+\dots+n) \\
 &= \frac{1}{n} \cdot \frac{n(n+1)}{2} = \frac{n+1}{2} \\
 \\
 \text{(ii) Variance } &= \sum_{i=1}^n x_i^2 f(x_i) - \bar{x}^2 \\
 &= 1^2 \cdot \frac{1}{n} + 2^2 \cdot \frac{1}{n} + 3^2 \cdot \frac{1}{n} + \dots + n^2 \cdot \frac{1}{n} - \\
 &\quad \left(\frac{n+1}{2}\right)^2 \\
 &= \frac{1}{n} (1^2 + 2^2 + 3^2 + \dots + n^2) - \frac{1}{4} (n+1)^2 \\
 &= \frac{1}{n} \cdot \frac{n(n+1)(2n+1)}{6} - \frac{1}{4} (n+1)^2 \\
 &= \frac{(n+1)(2n+1)}{6} - \frac{(n+1)^2}{4} \\
 &= \frac{(n+1)}{2} \left(\frac{2n+1}{3} - \frac{n+1}{2} \right) \\
 &= \frac{n+1}{12} (4n+2-3n-3) \\
 &= \frac{(n+1)(n-1)}{12} = \frac{n-1}{12}.
 \end{aligned}$$

Example: From a lot of 10 items containing 3 defectives, a sample of 4 items is drawn at random. Let the random variable X denote the number of defective items in the sample. Find the probability distribution of X when the sample is drawn without replacement.

Solution: X takes the values 0, 1, 2 or 3

Given total no. of items = 10

No. of good items = 7

No. of defective items = 3

$$P(X=0) = P(\text{No defective})$$

$$= \frac{7C_4}{10C_4} = \frac{7!}{4!3!} \times \frac{4!6!}{10!} = \frac{1}{6}$$

$$P(X=1) = P(\text{one defective and 3 good})$$

$$= \frac{3C_1 \times 7C_3}{10C_4} = \frac{3 \times 7!}{3!4!} \times \frac{4!6!}{10!} = \frac{1}{2}$$

$$P(X=2) = P(\text{2 defective and 2 good})$$

$$= \frac{3C_2 \times 7C_2}{10C_4} = \frac{3}{10}$$

$P(X=3) = P(3 \text{ defective and 1 good})$

$$= \frac{3C_3 \times 7C_1}{10C_4} = \frac{1}{10C_4} = \frac{4!}{8 \times 9 \times 10} = \frac{1}{30}$$

∴ The probability distribution of random variable X is as follows:

X	0	1	2	3
$P(X)$	$\frac{1}{6}$	$\frac{1}{2}$	$\frac{3}{10}$	$\frac{1}{30}$

Example: Find the distribution function which corresponds to the probability distribution defined by $f(x) = \frac{x}{15}$ for $x = 1, 2, 3, 4, 5$.

Solution:

Given $f(x) = \frac{x}{15}$

$$\Rightarrow f(1) = \frac{1}{15}, f(2) = \frac{2}{15}, f(3) = \frac{3}{15}$$

$$f(4) = \frac{4}{15}, f(5) = \frac{5}{15}$$

$$\text{Now, } F(1) = f(1) = \frac{1}{15},$$

$$F(2) = F(1) + f(2) = \frac{1}{15} + \frac{2}{15} = \frac{1}{5}$$

$$F(3) = F(2) + f(3) = \frac{1}{5} + \frac{1}{3} = \frac{8}{15}$$

$$F(4) = F(3) + f(4) = \frac{2}{5} + \frac{4}{15} = \frac{2}{3}$$

$$\text{and } F(5) = F(4) + f(5) = \frac{2}{3} + \frac{1}{3} = 1$$

2.10 General properties of a Continuous Probability Distribution:

(i) Mean:

Mean of a continuous probability distribution is given by

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

If X is defined from a to b ,

$$\text{Then, } \mu = \int_a^b x f(x) dx$$

(ii) Variance:

Variance of a continuous probability distribution is given by

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

(or)

$$\sigma^2 = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

If X is defined from a to b ,

$$\text{Then, } \sigma^2 = \int_a^b x^2 f(x) dx - \mu^2$$

(iii) Mode:

Mode is the value of x for which $f(x)$ is maximum. Thus mode is given by $f'(x)=0$ and $f''(x)<0$ for $a < x < b$.

(iv) Median:

Median is the point which divides the entire distribution into two equal parts. For continuous distribution, median is a point which divides the total area into two equal parts.

If x is defined a to b and M is the median, then

$$\int_a^M f(x) dx = \int_M^b f(x) dx = \frac{1}{2}$$

(v) Mean deviation:

Mean deviation about the mean (μ) is given by $\int_{-\infty}^{\infty} |x - \mu| f(x) dx$.

Example: For the continuous probability function $f(x) = Kx^2 e^{-x}$ when $x \geq 0$, find (i) K (ii) Mean (iii) Variance

Solution: (i) we know that $\int_{-\infty}^{\infty} f(x) dx = 1$

$$\Rightarrow \int_0^{\infty} Kx^2 e^{-x} dx = 1 \quad (\because x \geq 0)$$

$$\Rightarrow K[x^2(-e^{-x}) - 2x(e^{-x}) + 2(-e^{-x})]_0^{\infty} = 1$$

$$\Rightarrow K[-e^{-x}(x^2 + 2x + 2)]_0^{\infty} = 1$$

$$\Rightarrow K(0+2) = 1$$

$$\Rightarrow \boxed{K = \frac{1}{2}}$$

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

$$= \int_0^{\infty} Kx^3 e^{-x} dx$$

$$= K[x^3(-e^{-x}) - 3x^2(e^{-x}) + 6x(-e^{-x}) - 6(-e^{-x})]_0^{\infty}$$

$$= K \left[-e^{-x} (x^3 + 3x^2 + 6x + 6) \right]_0^\infty$$

$$= K [0 + 6] = 6K = 6\left(\frac{1}{2}\right) = 3$$

$$\therefore \boxed{\mu = 3}$$

$$(iii) \text{ Variance} = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

$$= \int_0^{\infty} x^2 \cdot K x^2 e^{-x} dx - 3^2$$

$$= K \int_0^{\infty} x^4 e^{-x} dx - 9$$

$$= K \left[x^4 (-e^{-x}) - 4x^3 (-e^{-x}) + 12x^2 (-e^{-x}) - 24x (-e^{-x}) + 24 e^{-x} \right]_0^\infty - 9$$

$$= \frac{1}{2} \left[-e^{-x} (x^4 + 4x^3 + 12x^2 + 24x + 24) \right]_0^\infty - 9$$

$$= \frac{1}{2} [0 + 24] - 9$$

$$\boxed{\therefore K = \frac{1}{2}}$$

$$= 12 - 9 = 3$$

Example: Suppose a continuous random variable X has the probability density $f(x) = K(1-x)$ for $0 < x < 1$, and $f(x) = 0$ otherwise. Find (i) K (ii) mean (iii) variance.

Solution : (i) we know that, $\int_{-\infty}^{\infty} f(x) dx = 1$

$$\text{i.e., } \int_{-\infty}^0 f(x) dx + \int_0^1 f(x) dx + \int_1^{\infty} f(x) dx = 1$$

$$\text{i.e., } 0 + \int_0^1 K(1-x)^2 dx + 0 = 1$$

$$\text{i.e., } K \left(x - \frac{x^3}{3} \right)_0^1 = 1$$

$$\text{i.e., } K \left(1 - \frac{1}{3} \right) = 1$$

$$\therefore \boxed{K = \frac{3}{2}}$$

$$(ii) \text{ Mean } \mu = \int_{-\infty}^{\infty} x f(x) dx$$

$$= \int_0^1 x \cdot f(x) dx$$

$$= \int_0^1 x \cdot K(1-x)^2 dx$$

$$\mu = K \int_0^1 x(x - \frac{x^3}{3}) dx$$

$$\Rightarrow \mu = K \left(\frac{x^2}{2} - \frac{x^4}{4} \right)_0^1 = K \left(\frac{1}{2} - \frac{1}{4} \right) = \frac{K}{4}$$

$$\mu = \left(\frac{3}{2} \right) \left(\frac{1}{4} \right) = \frac{3}{8} \quad [\because K = \frac{3}{2}]$$

$$\therefore \boxed{\mu = \frac{3}{8}}$$

$$\begin{aligned}
 \text{(iii) Variance}(\sigma^2) &= \int_0^1 x^2 f(x) dx - \bar{x}^2 \\
 &= \int_0^1 x^2 K(1-x)^4 dx - \left(\frac{3}{8}\right)^2 \\
 &= K \int_0^1 (x^2 - \bar{x}^2) dx - \frac{9}{64} \\
 \sigma^2 &= K \left(\frac{x^3}{3} - \frac{x^5}{5} \right) \Big|_0^1 - \frac{9}{64} \\
 &= K \left(\frac{1}{3} - \frac{1}{5} \right) - \frac{9}{64} \\
 &= \frac{2K}{15} - \frac{9}{64} \\
 &= \frac{2}{15} \left(\frac{3}{2} \right) - \frac{9}{64} \\
 &= \frac{1}{5} - \frac{9}{64} \\
 &= \frac{19}{320} = 0.06
 \end{aligned}$$

$\therefore \boxed{\sigma^2 = 0.06}$

Example: The Probability density function of a continuous random variable is given by $f(x) = C e^{-Cx}$, $-\infty < x < \infty$.
 Find (i) C (ii) mean (iii) variance (iv) $P(0 \leq X \leq 4)$

Solution: Given function $f(x) = C e^{-|x|}$, $-\infty < x < \infty$

(i) we know that $\int_{-\infty}^{\infty} f(x) dx = 1$

$$\Rightarrow \int_{-\infty}^{\infty} C e^{-|x|} dx = 1$$

$$\Rightarrow C \int_{-\infty}^{\infty} e^{-|x|} dx = 1$$

$$\Rightarrow 2C \int_0^{\infty} e^{-x} dx = 1 \quad [\because e^{-|x|} \text{ is even}]$$

$$\Rightarrow 2C \int_0^{\infty} e^{-x} dx = 1$$

$$\Rightarrow 2C (-e^{-x}) \Big|_0^{\infty} = 1 \quad [\because 0 \leq x \leq \infty, |x| = x]$$

$$\Rightarrow -2C (0 - 1) = 1$$

$$\Rightarrow 2C = 1$$

$$\Rightarrow \boxed{C = \frac{1}{2}}$$

$$\therefore f(x) = \frac{1}{2} e^{-|x|}$$

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

$$= \frac{1}{2} \int_{-\infty}^{\infty} x e^{-|x|} dx = 0 \quad [\because \text{Integrand is odd}]$$

$$(ii) \text{ Variance} = \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

$$= \int_{-\infty}^{\infty} (x - 0)^2 \frac{1}{2} e^{-|x|} dx$$

$$= \frac{1}{2} \int_{-\infty}^{\infty} x^2 e^{-|x|} dx$$

$$= 2 \cdot \frac{1}{2} \int_0^{\infty} x^2 e^{-x} dx$$

$\left[\because \text{Integrand is even} \right]$

$$= \int_0^{\infty} x^2 e^{-x} dx$$

$$= \left(x \frac{e^{-x}}{-1} - 2x \frac{e^{-x}}{-1} + 2 \frac{e^{-x}}{-1} \right)_0^{\infty}$$

$$= (0 - (-2)) = 2$$

$$(iii) P(0 \leq x \leq 4) = \frac{1}{2} \int_0^4 e^{-x} dx$$

$$= \frac{1}{2} \int_0^4 e^{-x} dx \quad \left[\because \text{In } 0 \leq x \leq 4, |x| = x \right]$$

$$= -\frac{1}{2} (\bar{e}^{-x})_0^4$$

$$= -\frac{1}{2} (\bar{e}^{-4} - 1)$$

$$= \frac{1}{2} (1 - \bar{e}^{-4}) = 0.4908$$

Example: A continuous random variable has the probability density function

$$f(x) = \begin{cases} Kx^{-\lambda x}, & \text{for } x \geq 0, \lambda > 0 \\ 0, & \text{otherwise.} \end{cases}$$

- (i) K
- (ii) mean
- (iii) variance

Solution: (i) we know that $\int_{-\infty}^{\infty} f(x) dx = 1$

$$\Rightarrow \int_{-\infty}^0 0 \cdot dx + \int_0^{\infty} Kx^{-\lambda x} dx = 1$$

$$\Rightarrow K \int_0^{\infty} x^{-\lambda x} dx = 1$$

$$\Rightarrow K \left[x \left(\frac{e^{-\lambda x}}{-\lambda} \right) - 1 \left(\frac{e^{-\lambda x}}{-\lambda} \right) \right]_0^{\infty} = 1$$

$$\Rightarrow K \left[(0 - 0) - (0 - \frac{1}{\lambda}) \right] = 1$$

$$\Rightarrow K = \lambda^{\nu}$$

$$\therefore f(x) = \begin{cases} \lambda^{\nu} x^{-\lambda x}, & \text{for } x \geq 0, \lambda > 0 \\ 0, & \text{otherwise} \end{cases}$$

av

$$(ii) \text{ Mean} = \mu = \int_{-\infty}^{\infty} x f(x) dx$$

$$\Rightarrow \mu = \int_{-\infty}^0 0 \cdot dx + \int_0^{\infty} x \lambda^x x e^{-\lambda x} dx$$

$$= \lambda \int_0^{\infty} x^2 e^{-\lambda x} dx$$

$$= \lambda \left[x^2 \left(\frac{-\lambda x}{(-\lambda)} \right) - 2x \left(\frac{e^{-\lambda x}}{\lambda^2} \right) + 2 \left(\frac{e^{-\lambda x}}{-\lambda^3} \right) \right]_0^\infty$$

$$= \lambda \left[(0 - 0 + 0) - (0 - 0 - \frac{2}{\lambda^3}) \right]$$

$$\boxed{\mu = \frac{2}{\lambda}}$$

$$(iii) \text{ Variance} = \sigma^2 = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

$$\Rightarrow \sigma^2 = \int_0^{\infty} x^2 f(x) dx - \left(\frac{2}{\lambda} \right)^2$$

$$= \lambda \int_0^{\infty} x^2 \lambda^x e^{-\lambda x} dx - \frac{4}{\lambda^2}$$

$$= \lambda \left[x^2 \left(\frac{-\lambda x}{(-\lambda)^2} \right) - 3x^2 \left(\frac{e^{-\lambda x}}{\lambda^2} \right) + 6x \left(\frac{e^{-\lambda x}}{-\lambda^3} \right) + 6 \left(\frac{e^{-\lambda x}}{\lambda^4} \right) \right]_0^\infty$$

$$= \lambda \left[(0 - 0 + 0 - 0) - (0 - 0 + 0 - \frac{6}{\lambda^4}) \right] - \frac{4}{\lambda^2}$$

$$\sigma^2 = \frac{6}{\lambda^2} - \frac{4}{\lambda^2} = \frac{2}{\lambda^2}$$

Example: A Continuous random variable X
has the distribution function

$$F(x) = \begin{cases} 0, & \text{if } x \leq 1 \\ K(x-1)^4, & \text{if } 1 < x \leq 3 \\ 1, & \text{if } x > 3 \end{cases}$$

Find (i) $f(x)$

(ii) K

(iii) Mean

Solution: (i) we know that $f(x) = \frac{d}{dx}[F(x)]$

$$\therefore f(x) = \begin{cases} 0, & \text{if } x \leq 1 \\ 4K(x-1)^3, & \text{if } 1 < x \leq 3 \\ 0, & \text{if } x > 3 \end{cases}$$

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

$$\Rightarrow \int_{-\infty}^0 f(x) dx + \int_1^3 f(x) dx + \int_3^{\infty} f(x) dx = 1$$

$$\Rightarrow 0 + \int_1^3 4K(x-1)^3 dx + 0 = 1$$

$$\Rightarrow 4K \left[\frac{(x-1)^4}{4} \right]_1^3 = 1$$

$$\Rightarrow K(16 - 0) = 1$$

$$\Rightarrow \boxed{K = \frac{1}{16}}$$

$$\therefore f(x) = \begin{cases} 0 & , \text{ if } x \leq 1 \\ \frac{1}{4}(x-1)^3 & , \text{ if } 1 < x \leq 3 \\ 0 & , \text{ if } x \geq 3 \end{cases}$$

$$\text{(iii) Mean} = \mu = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

$$= \int_{-\infty}^1 x \cdot f(x) dx + \int_1^3 x \cdot f(x) dx + \int_3^{\infty} x \cdot f(x) dx$$

$$= 0 + \int_1^3 x \cdot \frac{1}{4}(x-1)^3 dx + 0$$

$$= \frac{1}{4} \int_1^3 x(x-1)^3 dx$$

$$= \frac{1}{4} \int_0^2 (t+1)^3 t^3 dt \quad [\because \text{Put } x-1=t]$$

$$= \frac{1}{4} \int_0^2 (t^4 + t^3) dt$$

$$= \frac{1}{4} \left[\frac{t^5}{5} + \frac{t^4}{4} \right]_0^2$$

$$= \frac{1}{4} \left(\frac{2^5}{5} + \frac{2^4}{4} \right) = \frac{1}{4} \left(\frac{32}{5} + \frac{16}{4} \right)$$

$$= 4 \left(\frac{13}{20} \right) = \frac{13}{5} = 2.6$$

Example: If X is a continuous random variable and $Y = ax + b$, prove that $E(Y) = aE(X) + b$ and $V(Y) = a^2 V(X)$ where V is variance and a, b are constants.

Solution: (i) By definition,

$$\begin{aligned} E(Y) &= E(ax+b) = \int_{-\infty}^{\infty} (ax+b) f(x) dx \\ &= a \int_{-\infty}^{\infty} x f(x) dx + b \int_{-\infty}^{\infty} f(x) dx \\ &= aE(X) + b \quad [\because \text{Total probability is unity}] \\ E(Y) &= aE(X) + b \end{aligned}$$

$$(ii) \text{ we have } E(Y) = aE(X) + b \quad [\because \text{From (i)}]$$

$$\text{where } Y = ax + b \quad \textcircled{2}$$

$$\textcircled{2} - \textcircled{1} \Rightarrow Y - E(Y) = a[X - E(X)]$$

squaring on both sides

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

Taking expectation on both sides

$$\text{we get } E\{[Y - E(Y)]^2\} = a^2 E\{[X - E(X)]^2\}$$

$$\therefore V(Y) = a^2 V(X)$$

Example: If X is a continuous random variable and K is a constant, then prove that (i) $\text{Var}(X+K) = \text{Var}(X)$ (ii) $\text{Var}(KX) = K^2 \text{Var}(X)$

Solution:

By definition,

$$\text{Var}(X) = \int_{-\infty}^{\infty} x^2 f(x) dx - \left[\int_{-\infty}^{\infty} x f(x) dx \right]^2$$

$$\begin{aligned} \text{(i)} \quad \text{Var}(X+K) &= \int_{-\infty}^{\infty} (x+K)^2 f(x) dx - \left[\int_{-\infty}^{\infty} (x+K) f(x) dx \right]^2 \\ &= \int_{-\infty}^{\infty} (x^2 + 2Kx + K^2) f(x) dx - \left[\int_{-\infty}^{\infty} x f(x) dx + K \int_{-\infty}^{\infty} f(x) dx \right]^2 \\ &= \int_{-\infty}^{\infty} x^2 f(x) dx + 2K \int_{-\infty}^{\infty} x f(x) dx + K^2 - \\ &\quad \left[\int_{-\infty}^{\infty} x f(x) dx + K \right]^2 \\ &\quad \left[\because \int_{-\infty}^{\infty} f(x) dx = 1 \right] \\ &= E(X^2) + 2KE(X) + K^2 - [E(X) + K]^2 \end{aligned}$$

$$\begin{aligned}
 &= E(X^2) + 2KE(X) + K^2 - [E(X)]^2 - 2KE(X) - K^2 \\
 &= E(X^2) - [E(X)]^2 \\
 &= \text{var}(X)
 \end{aligned}$$

$$\begin{aligned}
 \text{(i)} \quad \text{var}(kX) &= \int_{-\infty}^{\infty} k^2 x^2 f(x) dx - \left[\int_{-\infty}^{\infty} kx f(x) dx \right]^2 \\
 &= k^2 \int_{-\infty}^{\infty} x^2 f(x) dx - k^2 \left[\int_{-\infty}^{\infty} f(x) dx \right]^2 \\
 &= k^2 [E(X)^2 - \{E(X)\}^2] \\
 &= k^2 \text{var}(X)
 \end{aligned}$$

2.11 Binomial Distribution

Binomial Distribution was discovered by James Bernoulli in the year 1700 and it is a Discrete Probability distribution.

Let us discuss a conceptual or particular situation where a trial or an experiment results in only two outcomes, say 'Success' and 'Failure'. Further, the result of one trial does not influence the result of next trial, and the probability of success at each trial is the same from trial to trial.

Some of such situations are:

1. Tossing a coin - Head or Tail
2. Birth of a baby - Girl or Boy.
3. Auditing a Bill - Contains an error or not.

The conditions for the applicability of a Binomial distribution are as follows,

- i) There are n independent trials
- ii) Each trial has only two possible outcomes.

iii) The probabilities of two outcomes remain constant.

→ Let the number of trials be n . The trials be independent i.e. the success or failure at one trial does not affect the outcome of the other trials. Thus the probability of success remains the same from trial to trial. Also let 'p' be the probability of success and 'q' be the probability of failure. Then we have $p+q=1$ or $q=1-p$.

The probability of getting such a sequence of r successes and $(n-r)$ failures is $= p^r q^{n-r}$ (using Multiplication Theorem on Probability)

However, the number of ways in which one can get r successes and $n-r$ failures is ${}^n C_r$, and the probability of each of these ways is $p^r q^{n-r}$.

If x is the random variable, representing the number of successes, the probability of getting r successes and $n-r$ failures, in 'n' trials, is given by the probability function

$$P(x=r) = {}^n C_r P^r q^{n-r}, \quad r=0, 1, 2, \dots, n.$$

The probability of the number of successes so obtained is called the Binomial Probability Distribution, because these probabilities are the successive terms in the expansion of the binomial $(q+p)^n$. In the expansion of the distribution, 'n' is known as the degree of the distribution.

Definition: A random variable x has a Binomial distribution if it assumes only non negative values and its probability density function is given by

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

→ The Binomial distribution function is given by $F_x(x) = P(X \leq x) = \sum_{r=0}^n nC_r p^r q^{n-r}$

Examples of Binomial Distribution:

- i) The number of defective bolts in a box containing 'n' bolts.
- ii) The number of post-graduates in a group of 'n' men.

Conditions of Binomial Distribution:

1. Trials are repeated under identical conditions for a fixed number of times, say n times.
2. There are only two possible outcomes, e.g. success or failure for each trial.
3. The probability of success in each trial remains constant and does not change from trial to trial.

4. The trials are independent i.e., the probability of an event in any trial is not affected by the results of any other trial.

2.11.1 Constants of Binomial Distribution:

① Mean of the Binomial Distribution:

The Binomial probability Distribution is

given by

$$P(r) = nCr p^r q^{n-r}; \quad r=0, 1, 2, \dots, n \text{ and} \\ q = 1-p.$$

$$\begin{aligned} \text{Mean of } X, \mu &= E(X) = \sum_{r=0}^n r P(r) \\ &= 0 \times q^n + 1 \times nC_1 p q^{n-1} + 2 \times p \\ &\quad + \dots + n \cdot p^n \\ &= npq^{n-1} + 2 \cdot \frac{n(n-1)}{2!} p^2 q^{n-2} + \end{aligned}$$

$$3 \frac{n(n-1)(n-2)}{3!} p^3 q^{n-3} + \dots + np^n.$$

$$= np(q+p)^{n-1} \quad (\text{using Binomial distribution})$$

$$= np(1) = np.$$

Hence Arithmetic mean of the Binomial distribution
= np.

(2) Variance of the Binomial Distribution.

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

$$\therefore \boxed{V(X) = npq}.$$

Hence the standard Deviation of the
Binomial Distribution = \sqrt{npq}

(3). Mode of the Binomial Distribution:

Mode of the binomial distribution is the
Value of x at which $p(x)$ has max. value.

$$\text{Mode} = \begin{cases} \text{integral part of } (n+1)p, & \text{if } (n+1)p \text{ is} \\ & \text{not an integer.} \\ (n+1)p \text{ and } (n+1)(p-1), & .. \end{cases}$$

Recurrence Relation for the Binomial
Distribution:

$$P(r+1) = \frac{(n-r)p}{(r+1)q} \cdot P(r)$$

Example: A fair coin is tossed six times.
Find the probability of getting four heads.

Solution: P = Probability of getting a head $= \frac{1}{2}$

q = probability of not getting head
 $= \frac{1}{2}$ and $n = 6$, $r = 4$

We know that $P(r) = nCr p^r q^{n-r}$

$$\therefore P(4) = 6C_4 \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^{6-4}$$

$$= \frac{6!}{4!2!} \left(\frac{1}{2}\right)^6$$

$$= \frac{6 \times 5}{2} \times \frac{1}{2^6} = \frac{15}{64} = 0.2344$$

Example: Determine the probability of getting the sum 6 exactly 3 times in 7 throws with a pair of fair dice.

Solution: In a single throw of a pair of fair dice, a sum of 6 can occur in 5 ways,

$(1,5), (5,1), (2,4), (4,2)$ and $(3,3)$ out of $6 \times 6 = 36$ ways,

Thus P = Probability of occurrence of 6

in one throw = $\frac{5}{36}$

$$q = 1 - P = 1 - \frac{5}{36} = \frac{31}{36}$$

n = Number of trials = 7

\therefore Probability of getting 6 exactly thrice in 7 throws

$$= 7C_3 P^3 q^{7-3}$$

$$= 7C_3 \left(\frac{5}{36}\right)^3 \left(\frac{31}{36}\right)^4$$

$$= \frac{35(125)(31)^4}{(36)^7} = 0.0516 \text{ (nearly)}$$

Example: Ten coins are thrown simultaneously
Find the probability of getting at least

- (i) Seven heads
- (ii) Six heads
- (iii) One head.

Solution? P = Probability of getting a head = $\frac{1}{2}$

q = Probability of not getting a head = $\frac{1}{2}$

The probability of getting x heads in a throw of 10 coins is

$$P(X=x) = P(x) = {}^{10}C_x \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{10-x};$$

$$x=0, 1, 2, \dots, 10$$

(i) Probability of getting at least seven heads is given by

$$P(X \geq 7) = P(X=7) + P(X=8) + P(X=9) +$$

$$P(X=10)$$

$$= {}^{10}C_7 \left(\frac{1}{2}\right)^7 \left(\frac{1}{2}\right)^3 + {}^{10}C_8 \left(\frac{1}{2}\right)^8 \left(\frac{1}{2}\right)^2 +$$

$${}^{10}C_9 \left(\frac{1}{2}\right)^9 \left(\frac{1}{2}\right) + {}^{10}C_{10} \left(\frac{1}{2}\right)^{10}.$$

$$= \frac{1}{2^{10}} \left[{}^{10}C_7 + {}^{10}C_8 + {}^{10}C_9 + {}^{10}C_{10} \right]$$

$$= \frac{1}{2^{10}} (120 + 45 + 10 + 1)$$

$$= \frac{176}{1024} = 0.1719.$$

(ii) This is left as an

$$P(X \geq 6) = \frac{1}{2^{10}} ({}^{10}C_6 + {}^{10}C_7 + {}^{10}C_8 + {}^{10}C_9 +$$

$$\begin{aligned}
 \text{(iii)} \quad P(\text{at least } 1 \text{ head}) &= P(Y \geq 1) \\
 &= 1 - P(Y=0) \\
 &= 1 - 10C_0 \left(\frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^{10} \\
 &= 1 - 1 \times \left(\frac{1}{2}\right)^{10} \\
 &= 1 - \frac{1}{2^{10}}.
 \end{aligned}$$

Example: The mean and Variance of a binomial distribution are 2 and $\frac{8}{5}$. Find n .

Solution: We are given,

$$\begin{aligned}
 \text{Mean of the Binomial distribution} &= 2 \text{ i.e} \\
 np &= 2 \quad \dots \textcircled{1}
 \end{aligned}$$

$$\text{Variance of the Binomial distribution} = \frac{8}{5} \text{ i.e}$$

$$npq = \frac{8}{5} \quad \dots \textcircled{2}$$

$$\frac{\textcircled{2}}{\textcircled{1}} \Rightarrow \frac{npq}{np} = \frac{\frac{8}{5}}{2} = \frac{4}{5} \Rightarrow q = \frac{4}{5}$$

$$p = 1 - q = 1 - \frac{4}{5} = \frac{1}{5}$$

$$n = \frac{2}{p} = \frac{2}{\frac{1}{5}} = 10 \quad \therefore n = 10.$$

Example:

In eight throws of a die 5 or 6 is considered a success. Find the mean number of success and the standard deviation.

Solution: P = The probability of

$$\text{success} = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}.$$

q = The probability of failure

$$= 1 - p = 1 - \frac{1}{3} = \frac{2}{3}.$$

n = Number of throws = 8

$$\therefore \text{Mean} = np = 8\left(\frac{1}{3}\right) = \frac{8}{3}$$

$$\text{Variance} = npq = (np)q = \left(\frac{8}{3}\right)\left(\frac{2}{3}\right)$$

$$= \frac{16}{9}$$

$$\text{Hence standard deviation} = \sqrt{\text{Variance}}$$

$$= \sqrt{16/9} = 4/3$$

Example: Fit a binomial distribution to the following frequency distribution

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

Here n = number of trials = 6 and

N = total frequency = $\sum f_i = 200$

$$\therefore \text{Mean} = \frac{\sum f_i x_i}{\sum f_i} = \frac{25+104+174+128+80+24}{200} = \frac{535}{200} = 2.675$$

Now mean of the binomial distribution
 $= np$

$$\text{i.e. } np = 6p = 2.675$$

$$\therefore p = \frac{2.675}{6} = 0.446$$

$$\text{and } q = 1-p = 1-0.446 = 0.554$$

Hence the binomial distribution to be fitted is given by

$$N(q+p)^n = 200(0.554+0.446)^6$$

$$= 200(0.554+0.446)^6$$

$$= 200 [(0.554)^6 + 6C_1 (0.554)^5 (0.446) +$$

$$+ 6C_2 (0.554)^4 (0.446)^2 + 6C_3 (0.554)^3 (0.446)^3 +$$

$$6C_4 (0.554)^2 (0.446)^4 + 6C_5 (0.554) (0.446)^5 +$$

$$+ 6C_6 (0.446)^6]$$

$$= 200 [0.02891 + 0.1396 + 0.2809 + 0.3016 \\ + 0.1821 + 0.05864 + 0.007866]$$

$$= 5.782 + 27.92 + 56.18 + 60.32 + 36.42 \\ + 11.728 + 1.5732$$

The expected frequencies can be rounded off to the nearest integer to get expected frequencies as whole numbers.

\therefore The successive terms in the expansion give the expected or theoretical frequencies

which are

x	0	1	2	3	4	5	6
f	13	25	52	58	32	16	4
Expected or Theoretical frequency	6	28	56	60	36	12	2

2.12

Poisson Distribution

S.D Poisson introduced poisson distribution as a rare distribution of rare events i.e the events whose probability of occurrence is very small but the number of trials which could lead to the occurrence of the event, are very large.

Definition of Poisson Distribution:

A random variable X is said to follow a Poisson distribution if it assumes only non-negative values and its probability density function is given by

$$P(x, \lambda) = P(X=x) = \begin{cases} \bar{e}^\lambda \lambda^x ; & x=0,1,2,\dots \\ 0 , & \text{otherwise} \end{cases}$$

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

Note: 1. It should be noted that

$$\begin{aligned}\sum_{x=0}^{\infty} P(X=x) &= \sum_{x=0}^{\infty} \frac{e^{-\lambda} \lambda^x}{x!} \\ &= e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} \cdot e^{\lambda} = 1.\end{aligned}$$

Hence equation ① is a probability function.

2. The Poisson Distribution function is

$$\begin{aligned}F_x(x) = P(X \leq x) &= \sum_{y=0}^x P(y) \\ &= e^{-\lambda} \sum_{y=0}^x \frac{\lambda^y}{y!}, \quad y=0, 1, 2, \dots\end{aligned}$$

Examples of Poisson Distribution:

- (i) The number of defective electric bulbs manufactured by a reputed company.
- (ii) The number of telephone calls per minute at a switch board.
- (iii) The number of particles emitted by a radio-active substance.
- (iv) The number of persons born blind per year in a large city.

2.12.1:

Constants of the Poisson Distribution

1. The Mean of the Poisson Distribution.

$$\begin{aligned}
 \text{Mean} = E(X) &= \sum_{x=0}^{\infty} x \cdot P(x) = \sum_{x=0}^{\infty} x \cdot \frac{e^{-\lambda} \lambda^x}{x!} \\
 &= e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} \quad [\because x! = x(x-1)!] \\
 &= e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^{y+1}}{y!} \quad (\text{putting } y=x-1) \\
 &= \lambda e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} = e^{-\lambda} \cdot \lambda \cdot e^{\lambda} \\
 &\quad [\because \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} = e^{\lambda}] \\
 &= \lambda (= np)
 \end{aligned}$$

Thus the parameter λ is the Arithmetic Mean of the Poisson Distribution.

2. Variance of Poisson Distribution:

$$\begin{aligned}
 V(X) &= E(X^2) - [E(X)]^2 \\
 &= \sum_{x=0}^{\infty} x^2 P(x) - \lambda^2 \quad [\because \lambda = \text{mean of P.D}]
 \end{aligned}$$

$$= \sum_{x=0}^{\infty} x^2 \cdot \frac{e^{-\lambda} \lambda^x}{x!} - \lambda^2$$

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

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

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

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

$$= e^{-\lambda} \left[\sum_{y=0}^{\infty} \frac{\lambda^{y+2}}{y!} + \sum_{z=0}^{\infty} \frac{\lambda^{z+1}}{z!} \right] - \lambda^2$$

(putting $y = x-2$,
 $z = x-1$)

$$= e^{-\lambda} \left[\lambda^2 \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} + \lambda \sum_{z=0}^{\infty} \frac{\lambda^z}{z!} \right] - \lambda^2$$

$$\begin{aligned}
 &= e^{-1} [1^2 \cdot e^1 + 1 \cdot e^1] - 1^2 \\
 &= (1^2 + 1) - 1^2 \\
 &= 1.
 \end{aligned}$$

Thus Variance = 1

Hence the variance of the distribution

$$\text{Mean of the distribution} = 1$$

Further, standard deviation of the

$$\text{Poisson Distribution}, \sigma = \sqrt{\lambda}$$

3. Mode of the Poisson Distribution:

Mode is the value of x for which the

Probability $P(x)$ is maximum.

$$\therefore P(x) \geq P(x+1) \text{ and } P(x) \geq P(x-1)$$

$$\text{Now } P(x) \geq P(x+1) \Rightarrow \frac{e^{-1} \cdot 1^x}{x!} \geq \frac{e^{-1} \cdot 1^{x+1}}{(x+1)!}$$

$$\Rightarrow 1 \geq \frac{1}{x+1} \text{ or } \frac{1}{x+1} \leq 1$$

$$\Rightarrow 1 \leq x+1 \text{ or } x+1 \geq 1.$$

$$\Rightarrow x \geq 1-1$$

Similarly $P(x) \geq P(x-1)$

$$\Rightarrow x \leq \lambda$$

Combining (1) and (2), we have

$$\lambda - 1 \leq x \leq \lambda$$

Hence mode of the Poisson distribution

lies between $\lambda - 1$ and λ .

Case (1): If λ is an integer then $\lambda - 1$ is also an integer. So we have two maximum values and the distribution is bimodal and two modes are $(\lambda - 1)$ and λ .

Case (2): If λ is not an integer, the mode of Poisson distribution is integral part of λ .

4) Recurrence Relation for the Poisson Distribution:

$$\text{we have } P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$P(x+1) = \frac{e^{\lambda} \cdot \lambda^{x+1}}{(x+1)!} = \frac{1}{x+1} \cdot \frac{e^{\lambda} \lambda^x}{x!} = \frac{1}{x+1} P(x)$$

$$\text{Thus } P(x+1) = \left(\frac{1}{x+1}\right) P(x)$$

$$(OR) \quad P(x) = \frac{1}{x} \cdot P(x-1)$$

which is the required relation. which
this formula we can find $P(1), P(2), P(3), \dots$
if $P(0)$ is given.

Note: 1. on successive application of

the formula, we get

$$P(x+1) = \frac{\lambda^{x+1}}{(x+1)!}, \quad P(0)$$

$$2. \quad P(x \leq t) = \left[1 + \frac{1}{1!} + \frac{1^2}{2!} + \dots + \frac{1^t}{t!} \right] \cdot P(0)$$

$$3. \quad P(a \leq x \leq b) \left[\frac{1^a}{a!} + \frac{1^{a+1}}{(a+1)!} + \dots + \frac{1^b}{b!} \right] P(0)$$

for all $0 \leq a, b, a, b \in \mathbb{Z}^+$.

2.12.2

Properties of Poisson Distribution

1. Range of the variable is from 0 to ∞ .
2. Mean and Variance are equal.
3. Distribution gets more and more symmetrical about the mean as λ increases and tends to normal distribution, described in the next section.

Example: A hospital switch board receives an average of 4 emergency calls in a 10 min interval. What is the probability that (i) There are at most 2 emergency calls in a 10 minute interval (ii) There are exactly 3 emergency calls in a 10 minute interval.

$$\text{Mean, } \lambda = (\text{4 calls / 10 minutes}) = 4 \text{ calls}$$

$$\therefore P(X=x) = p(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$= \frac{e^{-4} 4^x}{x!} = \frac{1}{e^4} \cdot \frac{4^x}{x!}$$

$$(i) P(\text{at most 2 calls}) = P(X \leq 2)$$

$$+ P(X=2) = P(X=0) + P(X=1)$$

$$= \frac{1}{e^4} + \frac{1}{e^4} \cdot 4 + \frac{1}{e^4} \cdot \frac{4^2}{2!}$$

$$= \frac{1}{e^4} (1+4+8) = 13e^{-4} = 0.2381$$

(ii) $P(\text{exactly 3 calls}) = P(X=3)$

$$= \frac{1}{e^4} \cdot \frac{4^2}{3!} = \frac{32}{3} e^{-4}$$

$$= 0.1954$$

Example: If a bank received on the average 6 bad cheques per day, find the probability that it will receive 4 bad cheques on the same day.

Solution: we have $P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$

Here $\lambda=6$

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

Example: A manufacturer knows that the condensers he makes contain on average 1% defectives. He packs them in boxes of 100. What is the probability that a box

picked at random will contain 3 or more faulty condensers?

Solution: Here p = Probability of defective condensers = $1\% = 0.01$

n = Total number of condensers
 $= 100$

$$\therefore \text{Mean} = np = 100(0.01) = 1.$$

$$P(X=x) = P(x) = \frac{e^{-1} \cdot 1^x}{x!} = \frac{e^{-1} \cdot 1^x}{x!}$$
$$= \frac{\cancel{e^{-1}}}{\cancel{x!}}$$

$$P(X \geq 3) = 1 - P(X < 3)$$

$$= 1 - [P(X=0) + P(X=1) + P(X=2)]$$

$$= 1 - [e^{-1} + e^{-1} + \frac{e^{-1}}{2}]$$

$$= 1 - \cancel{e^{-1}} \cdot \frac{5}{2}$$

$$= 1 - 0.9197$$

$$= 0.0803.$$

Example: A manufacturer of cotter pins knows that 5% of his product is defective. Pins are sold in boxes of 100. He guarantees that not more than 10 pins will be defective. What is the approximate probability that a box will fail to meet the guaranteed quality?

Solution: The probability of cotter pins to be defective $= p = 5\% = 0.05$

Total number of cotter pins, $n=100$

$$\therefore \text{Mean, } \lambda = np = 100(0.05) = 5$$

$$\text{we have } P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\therefore P(X=x) = \frac{e^{-5} 5^x}{x!}$$

$P(\text{a box will fail to meet the guarantee})$

$$= P(X > 10) = 1 - P(X \leq 10)$$

$$= 1 - [P(X=0) + P(X=1) + \dots + P(X=10)]$$

$$= 1 - \left[\frac{e^{-5}(5)^0}{0!} + \frac{e^{-5}(5)^1}{1!} + \frac{e^{-5}(5)^2}{2!} + \frac{e^{-5}(5)^3}{3!} + \dots + \frac{e^{-5}(5)^{10}}{10!} \right] = 1 - 0.9863 \\ = 0.0137$$

Example: If a random variable has a Poisson distribution such that $P(1) = P(2)$.

Find (i) mean of the distribution

$$(ii) P(4) \quad (iii) P(x \geq 1) \quad (iv) P(1 < x < 4)$$

Solution: Given $P(1) = P(2)$

$$\text{Hence } \frac{e^{-\lambda} \lambda^1}{1!} = \frac{e^{-\lambda} \lambda^2}{2!}$$

$$\Rightarrow \lambda^2 = 2\lambda$$

$$\Rightarrow \lambda^2 - 2\lambda = 0 \Rightarrow \lambda(\lambda - 2) = 0$$

$$\therefore \lambda = 0 \text{ or } 2$$

(i) But $\lambda \neq 0$, $\therefore \lambda = 2$ Hence mean of the distribution = 2

$$(ii) P(x) = \frac{e^{-\lambda} \lambda^x}{x!} = \frac{e^{-2} \cdot 2^x}{x!}$$

$$P(x=4) = P(4) = \frac{e^{-2} \cdot 2^4}{4!} = \frac{2}{3e^2} \\ = 0.09022$$

$$(iii) P(x \geq 1) = 1 - P(x < 1) = 1 - P(x=0)$$

$$= 1 - \frac{e^{-2} \cdot 2^0}{0!} = 1 - \frac{1}{e^2} = 0.8647.$$

$$(iv) P(1 < x < 4) = P(x=2) + P(x=3)$$

$$= \frac{e^{-2} \cdot 2^2}{2!} + \frac{e^{-2} \cdot 2^3}{3!} = 2e^{-2} + \frac{4}{3}e^{-2} \\ = 0.4511$$

Example: Average number of accidents on any day on a national highway is 1.8,

Determine the probability that the number of accidents are (i) at least one (ii) atmost one

Solution: Mean, $\lambda = 1.8$

$$\text{we have } P(X=x) = P(x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

$$= \frac{e^{-1.8} \cdot (1.8)^x}{x!}$$

$$\begin{aligned} \text{(i)} \quad P(\text{at least one}) &= P(X \geq 1) \\ &= 1 - P(X=0) \\ &= 1 - e^{-1.8} \\ &= 1 - 0.1653 \\ &= 0.8347. \end{aligned}$$

$$\begin{aligned} \text{(ii)} \quad P(\text{at most one}) &= P(X \leq 1) \\ &= P(X=0) + P(X=1) \\ &= e^{-1.8} + e^{-1.8} (1.8) \\ &= e^{-1.8} (2.8) \\ &= 0.4628. \end{aligned}$$

Example: Number of monthly breakdowns of a computer is a random variable having poisson distribution with mean equal to 1.8.

Find the probability that the computer will function for a month
(i) without a breakdown
(ii) with only one breakdown and
(iii) with at least one

Solution: Given mean, $\lambda = 1.8$

from poisson distribution

$$P(X=x) = P(x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$
$$= \frac{(e^{-1.8})(1.8)^x}{x!}$$

$$(i) P(X=0) = \frac{(e^{-1.8})(1.8)^0}{0!} = e^{-1.8} = 0.1653$$

$$(ii) P(X=1) = \frac{(e^{-1.8})(1.8)^1}{1!} = e^{-1.8}(1.8) = 0.2975$$

$$(iii) P(X \geq 1) = 1 - P(X < 1)$$
$$= 1 - P(X=0)$$
$$= 1 - 0.1653$$
$$= 0.8347$$

Example: A distribution of bean seeds determines from extensive tests that 5% of large batch of seeds will not germinate. He sells the seeds in packets of 200 and guarantees 90% germination. Determine the probability that a particular packet will violate the guarantee.

Solution: P = The probability of a seed

not germinating = 5% = 0.05

λ = mean number of seeds in a

sample of 200 = $\lambda = np = 200 \times 0.05 = 10$

Let x be the number of seeds that do not germinate. Then

$$P(x) = \frac{e^{-10} \cdot 10^x}{x!}$$

A packet will violate guarantee if it contains more than 20 germination seeds.

Probability that the guarantee is violated = $P(x > 20)$

$$= 1 - P(x \leq 20)$$

= e^{-10}

$$= 1 - \sum_{x=0}^{20} \frac{e^{-10} 10^x}{x!}$$

$$= 1 - 0.9984$$

$$= 0.0016$$

Example: If a Poisson distribution is

such that $P(x=1) \cdot \frac{3}{2} = P(x=3)$,

find (i) $P(x \geq 1)$ (ii) $P(x \leq 3)$ (iii) $P(2 \leq x \leq 5)$

Solution: Given $\frac{3}{2} P(x=1) = P(x=3)$

$$\text{i.e. } \frac{3}{2} \cdot \frac{e^{-\lambda} \lambda^1}{1!} = \frac{e^{-\lambda} \lambda^3}{3!}$$

$$\text{i.e. } \frac{3\lambda}{2} = \frac{\lambda^3}{6}$$

$$\text{i.e. } \lambda^3 = 9\lambda \Rightarrow \lambda^3 - 9\lambda = 0$$

$$\Rightarrow \lambda(\lambda^2 - 9) = 0$$

$$(\text{or}) \quad \lambda = 0, \quad \lambda^2 = 9$$

$$\lambda = 0, \quad \lambda = 3, -3.$$

$$\Rightarrow \underline{\lambda = 3} \quad (\because \lambda > 0)$$

$$\text{Hence } P(x=x) = P(x) = \frac{e^{-3} 3^x}{x!}$$

$$\begin{aligned}
 \text{(i)} \quad P(X \geq 1) &= 1 - P(X < 1) \\
 &= 1 - P(X = 0) \\
 &= 1 - \frac{\bar{e}^3 \cdot 3^0}{0!} \\
 &= 1 - \bar{e}^3 \\
 &= 0.950213.
 \end{aligned}$$

$$\begin{aligned}
 \text{(ii)} \quad P(X \leq 3) &= P(X = 0) + P(X = 1) + P(X = 2) \\
 &\quad + P(X = 3) \\
 &= \bar{e}^3 \left[\frac{3^0}{0!} + \frac{3^1}{1!} + \frac{3^2}{2!} + \frac{3^3}{3!} \right] \\
 &= \bar{e}^3 \left(1 + 3 + \frac{9}{2} + \frac{9}{2} \right) \\
 &= 13 \bar{e}^3 \\
 &= 0.6472318.
 \end{aligned}$$

$$\begin{aligned}
 \text{(iii)} \quad P(2 \leq X \leq 5) &= P(X = 2) + P(X = 3) + \\
 &\quad P(X = 4) + P(X = 5) \\
 &= \bar{e}^3 \left[\frac{3^2}{2!} + \frac{3^3}{3!} + \frac{3^4}{4!} + \frac{3^5}{5!} \right] \\
 &= 9 \bar{e}^3 \left(\frac{1}{2} + \frac{1}{2} + \frac{3}{8} + \frac{9}{40} \right) \\
 &= 9 \bar{e}^3 (1.6) \\
 &= 0.7169337.
 \end{aligned}$$

Example: If the Variance of a Poisson Variate is 3, then find the probability that (i) $x=0$ (ii) $0 < x \leq 3$ (iii) $1 \leq x < 4$.

Solution: Given Variance of the Poisson distribution, we have

$$\text{mean} = \text{Variance}$$

Hence mean of the Poisson Distribution = 3,

$$\text{i.e } \lambda = 3,$$

$$\text{So } P(X=x) = P(x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

$$= \frac{e^{-3} \cdot 3^x}{x!}$$

$$(i) P(x=0) = \frac{e^{-3} \cdot 3^0}{0!} = e^{-3}$$

$$= 0.04979$$

$$(ii) P(0 < x \leq 3) = P(x=1) + P(x=2) + P(x=3)$$

$$= e^{-3} \cdot 3 + e^{-3} \cdot \frac{3^2}{2!} + e^{-3} \cdot \frac{3^3}{3!}$$

$$= e^{-3} \left(3 + \frac{9}{2} + \frac{27}{8} \right)$$

$$= e^{-3} (12) = 0.5974$$

$$\text{(iii)} \quad P(1 \leq x < 4) = P(x=1) + P(x=2) + P(x=3)$$
$$= \underline{0.5974}$$