

25/08/23

(UNIT-2)

Correlation and Regression

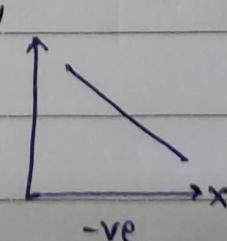
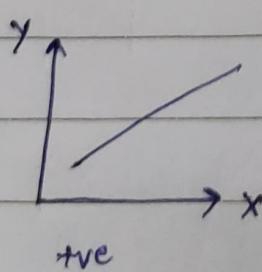
Topics

- ⇒ Scatter Plots.
- ⇒ Correlation coefficient and properties.
- ⇒ Karl Pearson's.
- ⇒ Spearman's Rank Correlation.
- ⇒ Regression.
- ⇒ Linear Regression and its Properties.

If X and Y are two RV and when they are inter-related.

$X \rightarrow$ Demand ↑ ↓	$\Rightarrow X \rightarrow$ Placement ↓	$\Rightarrow X \rightarrow$ attendance
$Y \rightarrow$ Supply ↑ ↓ ↓ +ve correlation	$Y \rightarrow$ Admission ↑ ↓ -ve correlation	$Y \rightarrow$ CGPA ↓ No correlation

- * If the change in one variable affect the changes of another variable, then the variables are correlated.
- * If the two variables deviate in the same direction, then the correlation is said to be direct or (+ve).
- * If two RV deviate in the opposite direction then the correlation is said to be diverse or negative (-ve).



* Scatter diagram or Plots.

↳ for the bivariate distribution (x_i, y_i) ($i=1, 2, \dots, n$), if the values of x_i and y_i are plotted along x axis and y -axis in the xy plane, the diagram of dots so obtained is known as scatter diagram/plots.

a) \Rightarrow If the points are very dense i.e., very close to each other, then x and y are correlated. (fairly good data).

b) \Rightarrow If the points are widely scattered, either a poor correlation or no correlation is expected.

\Rightarrow dense scatter

\checkmark (+ve) correlation

widely scatter

\checkmark (-ve) correlation

No correlation.

* Karl Pearson's Coefficient of Correlation.

Correlation coefficients b/w two RV X and Y are denoted by

$\gamma(x, y)$ or γ_{xy} which is defined as:-

$$\star \boxed{\gamma(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}} \quad \text{eq ①.}$$

$$\text{Cov}(x, y) = \frac{1}{m} \sum (x_i - \bar{x})(y_i - \bar{y})$$

means
no of elements

$$\bar{x}^2 = E\left\{ x - E(x) \right\}^2 = \frac{1}{m} \sum (x_i - \bar{x})^2$$

$$\bar{y}^2 = E\left\{ y - E(y) \right\}^2 = \frac{1}{m} \sum (y_i - \bar{y})^2$$

from eq ①

$$\begin{aligned} \gamma(x,y) = \gamma_{xy} &= \frac{\frac{1}{n} \sum xy - \bar{x}\bar{y}}{\sqrt{\frac{1}{n} \sum x^2 - \bar{x}^2} \cdot \sqrt{\frac{1}{n} \sum y^2 - \bar{y}^2}} \\ &= \frac{\sum xy - \bar{x}\bar{y}}{\sqrt{\sum x^2 - \bar{x}^2} \cdot \sqrt{\sum y^2 - \bar{y}^2}} \end{aligned}$$

- Q Calculate the correlation of height of father's and their sons in inches. fathers $\rightarrow x = 65, 66, 67, 68, 69, 70, 72$
 sons $\rightarrow y = 67, 68, 65, 68, 72, 72, 69, 71$

\rightarrow	x	y	$U = x - \bar{x}$ $= x - 68$	$V = y - \bar{y}$ $= y - 69$	U^2	V^2	UV
	65	67	-3	-2	9	164	126
	66	68	-2	-1	4	91	62
	67	65	-1	-4	1	416	24
	67	68	-1	-1	1	41	21
	68	72	0	3	0	19	0
	69	72	1	3	1	89	23
	70	69	2	0	4	10	20
	72	71	4	2	16	84	128
					36	44	3624

* we will assume any random value as a mean
 instead we assumed mean

We will not use it

$$\bar{x} = \frac{\sum x}{n} = \frac{544}{8} = 68.$$

$$\text{Cov} = \frac{24}{8} = 3$$

$$\bar{x}^2 = 4.5^2 = \frac{36}{8} = 4.5$$

$$\bar{y} = \frac{\sum y}{n} = \frac{552}{8} = 69.$$

$$\bar{y}^2 = 4.5^2 = \frac{20}{8} = 5.5$$

$$\gamma_{xy} = \frac{3}{\sqrt{4.5} \times \sqrt{5.5}}$$

$$r_{xy} = \frac{36}{18} \times \frac{8}{\sqrt{2} \times \sqrt{44}} = \frac{3}{\frac{6^2}{\sqrt{8}} \times \frac{\sqrt{44}}{\sqrt{8}}} = \frac{1}{2\sqrt{44}} = \frac{1}{2\sqrt{44}} = \frac{1}{\sqrt{44}} = \frac{1}{\sqrt{44}} = \frac{1}{\sqrt{44}} = \frac{1}{\sqrt{44}} = \frac{1}{\sqrt{44}}$$

* also when elements are factor of

5
10
15
25

$$U = \frac{x-a}{h}, V = \frac{y-b}{k}$$

a > mean origin

k > scale

31/08/23

Q A computer while calculating correlation coefficient b/w two variables x and y from 25 pairs of observation, obtained the following results.

$$n=25, \sum x = 125, \sum x^2 = 650, \sum y = 100, \sum y^2 = 460, \sum xy = 508,$$

it was however discovered at the time of checking that he had copied down two pairs as while the correct values are :-

x	y
8	12
6	8

x	y
6	14

obtain correct correlation coefficient.

$$\rightarrow r_{xy} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} = \frac{\sum xy - \bar{x}\bar{y}}{\sqrt{\sum x^2 - \bar{x}^2} \sqrt{\sum y^2 - \bar{y}^2}} \quad \text{--- (1)}$$

$$\sum x = 125 - 6 - 8 + 8 + 6 = 125$$

$$\sum xy = 508 - 96 - 58 + 84 + 48$$

$$\sum y = 100 + 12 + 8 - 14 - 6 = 100$$

$$= 496$$

$$\sum x^2 = 650 - 36 - 64 + 64 + 36 = 650$$

$$= 508 - 84 - 48 + 96 + 48$$

$$\sum y^2 = 460 - 196 - 36 + 144 + 64 = 4864$$

$$= 520$$

$$\downarrow \quad \downarrow \quad 436$$

$$\bar{x} = \frac{\sum x}{n} = \frac{125}{25} = 5 \quad \bar{y} = \frac{\sum y}{n} = \frac{100}{25} = 4,$$

$$= \frac{1}{25} (520) - 25$$

$\begin{array}{r} 25 \\ \times 16 \\ \hline 120 \\ + 25 \\ \hline 380 \end{array}$

$$\sqrt{125} (650) - 25 \quad \sqrt{125} (436) - 20$$

$$\begin{array}{r} \cancel{520} - \cancel{25} = \cancel{495} \\ \cancel{42} \qquad \qquad \qquad \cancel{46} \end{array}$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (\bar{x}_i - \bar{x})(\bar{y}_i - \bar{y})$$

$$= \frac{1}{25} (520) - 20$$

$$= 20.8 - 20 \\ = 0.8$$

$$r_{xy} = \frac{0.8 \times 5}{6} = \frac{4.8}{6} = 0.8$$

25
12 5
25 X
12 5

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

$$= \sqrt{\frac{1}{25}(650)} - 25$$

$$C_4 = \sqrt{\frac{1}{n_{425}} (426) - 16} \quad 425$$

$$2 \sqrt{436 - \cancel{500}^{400}} \frac{144}{25} \sqrt{\frac{36}{25}} = 6/\cancel{5}$$

$$2 \sqrt{\frac{25}{25}} = 1$$

* Theorem :-

The Variable X and Y are connected by the equation $ax + by + c = 0$. Show that the correlation coeff b/w them is -1, if the signs of 'a' & 'b' are alike (same) and +1, if the sign of 'a' and 'b' are different.

$$\Rightarrow \det ax+by+c = 0 \quad \text{--- ①}$$

Taking Expectation on both side.

$$aE(x) + bE(y) + c = 0 \quad -\textcircled{2}.$$

① - ②

$$\text{Cov}(x, y) = ?$$

$$\sigma_x^2 = ?$$

$$a(x - E(x)) + b(y - E(y)) + c - c = 0$$

* Properties of Correlation :-

⇒ Correlation coefficient always lies b/w "-1 to +1;"

$$\text{i.e. } [-1 \leq r \leq 1]$$

a) If $r = +1$, then correlation is perfect and positive (+ve)

b) If $r = -1$, then correlation is perfect and negative (-ve)



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

$$\text{i.e. } U = \frac{x-a}{h}, \quad V = \frac{y-b}{k}, \quad a, b \rightarrow \text{origin (mean assumed)} \\ k, h \rightarrow \text{scale. (factor)}$$

⇒ Two independent variables are un-correlated.

Proof :- If two rv's are independent, then

$$\text{cov}(x, y) = 0$$

$$\therefore r(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{0}{\sigma_x \cdot \sigma_y}$$

$$\boxed{r_{(x,y)} = 0}$$

* Rank Correlation :-

⇒ Spearman's rank correlation Coefficient :-

⇒ The rank corr. is denoted by 'ρ'.

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

where $d_i = x_i - y_i$

{when rank is not repeated}

1/09/23

Q The rank of ^{same} 16 student in maths and Physics are as:-
Two numbers with bracket, denote the rank of the student in
math and physics.

$\left\{ (1,1), (2,10), (3,3), (4,4), (5,5), (6,7), (7,2), (8,6) \right. \\ \left. (9,8), (10,11), (11,15), (12,9), (13,12), (14,12), (15,16), (16,13) \right\}$

calculate the rank corr coeff for colfficiency of this group for Maths.
and Physics.

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

$$= 1 - \frac{6(136)}{16(255)} = 1 - \frac{84}{364} = 0.607$$

$$r_d = 1 - \frac{1}{24.285}$$

$$= 1 - 0.2$$

$$= 0.8$$

X	Y	$d_i = x - y$	d_i^2
1	15	-14	0
2	10	-8	64
3	3	-2	0
4	5	-1	0
5	8	-3	0
6	7	-1	1
7	2	-5	25
8	6	-2	4
9	8	-1	1
10	11	-1	1
11	15	-4	16
12	9	3	9
13	14	-1	1
14	12	2	4
15	16	-1	1
16	13	2	4
			136

NOTE

⇒ Rank can be (ve) also.

— / —

Q 10 competitor in a musical test were ranked by three judges (A, B, C) in following order.

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

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

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

using rank correlation matter, discuss which pair of judges has the nearest approach to common liking in music.

\Rightarrow	X	Y	Z	$d_1 = X - Y$	$d_2 = X - Z$	$d_3 = Y - Z$	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	16
3	7	1		-4	2	6	16	4	36
42	10	2		-8	0	8	64	0	64
4	2	3		2	1	-1	4	1	1
9	1	10		8	-1	-9	64	1	81
7	6	5		1	2	1	1	4	1
8	9	7		-1	-1	2	1	1	4
							200	60	214

$$f(X, Y) = \frac{1 - \left(6 \sum d_i^2 \right)}{10(100-1)} \Rightarrow 1 - \frac{6 \times 204}{10 \times 9533} = 1 - \frac{120}{9533} = \frac{-71}{9533} = \boxed{-\frac{71}{9533}}$$

$$f(X, Z) = 1 - \frac{6 \times 214}{10 \times 9533} = 1 - \frac{36 + 214}{9533} = \boxed{\frac{7}{11}}$$

$$f(Y, Z) = 1 - \frac{6 \times 214}{5 \times 10 \times 9533} = 1 - \frac{214}{165} = \boxed{-\frac{49}{165}}$$

(B)

$\Rightarrow f(x, z)$ having the max. value, so
the judges A and C have common liking in music.

* Repeleted Rank :-

$$\Rightarrow P = 1 - 6 \left[\frac{\sum d^2 + m_1(m_1^2 - 1)}{12} + \frac{m_2(m_2^2 - 1)}{12} + \dots \right] \div n(n^2 - 1)$$

.. where $m_1, m_2 \dots$ no. of times rank is repeated.

Q. Obtain the rank co-relation for the following data :-

where data are:-

$$\Rightarrow X: 68, 64, 75, 50, 64, 80, 75, 40, 55, 64$$

$$Y: 62, 58, 68, 45, 81, 60, 68, 48, 50, 70$$

75
2 3
write here
in ascending

X	Y	Rank of (X)	Rank of (Y)	$d = x - y$	d^2	$\frac{75}{2} = 2+3 = 2.5$
68	62	4	5.5	-1	1	$64 \Rightarrow \frac{5+6+7}{3} = 18/3 = 6$
64	58	6	7	-1	1	
75	68	2.5	3.5	-1	1	$68 = \frac{3+4}{2} = 3.5$
50	45	9	10	-1	1	
64	81	6	1	5	25	
80	60	1	6	-5	25	
75	68	2.5	3.5	-1	1	
40	48	10	9	1	1	
55	50	8	8	0	0	
64	70	6	2	4	16	$\rightarrow 72$

$$\frac{1 - 6 \left[\sum d^2 + \frac{m_1(m_1^2 - 1)}{12} + \frac{m_2(m_2^2 - 1)}{12} + \dots \right]}{n(n^2 - 1)}$$

In X series

↳ 75 is repeated 2 times

↳ 64 is repeated 3 times.

$$\Rightarrow \cancel{1 - 6 \left[72 + \cancel{\dots} \right]}$$

$$\frac{m_1(m_1^2 - 1)}{12} = \frac{2(4-1)}{12} = \frac{5}{12}$$

$$\frac{m_2(m_2^2 - 1)}{12} = \frac{3(9-1)}{12} = \frac{24}{12}$$

$$= \frac{5}{12} + \frac{24}{12} = \frac{29}{12}$$

$$= \frac{5}{12}$$

$$\Rightarrow f = \frac{1 - 6 \left[72 + \frac{5}{12} + \frac{1}{2} \right]}{10(99)} = 1 - \frac{\frac{2}{6} \times 75}{\frac{10}{2} \times \frac{29}{33}}^{15}$$

$$= 1 - \frac{15}{33}$$

$$= \frac{18}{33}$$

$$= \textcircled{0.545} \quad \text{Ans}$$

In Y-series

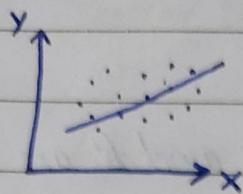
↳ 68 is repeated 2 times.

$$= \frac{m_2(m_2^2 - 1)}{12}$$

$$= \frac{2(4-1)}{12} = \frac{1}{2}$$

4/09/23

* Regression



* \rightarrow Best fit data

* \rightarrow Give future value and estimated.

where $\rightarrow x$ depends on y

\Rightarrow Regression analysis is a mathematical measure of the average relationship b/w two or more variables in terms of the original units of data.

$$x = a + by$$

constant
 depended var
 independent var
 slope

* Types of Regression :-

1) Linear Regression

2) Curvi - Linear Regression

(any line except straight line.)

* Linear Regression.

\Rightarrow If the variables of bi-variate distribution are related, then the points in the scattered diagram will cluster round some curve called the curve's of Regression.

\Rightarrow If the curve is a straight line it is called Line of regression, and said to be Linear regression b/w the variables, otherwise it will be curvi-linear regression.

\Rightarrow The line of regression is the line which gives the best estimate to the value of one variable for any specific value of other variable. Thus the line of regression is the line of "best fit".

Let us suppose that in bi-variate distribution, such that Y is dependent variable and X is independent variable.

$$(x_i, y_i) \text{ } i=1, 2, 3, \dots, n$$

Y on X is $\Rightarrow Y = a + bx$, where a and b are arbitrary const;
 $b \rightarrow$ is the scope of line of Regression.

The line of Regression passes through the points (\bar{x}, \bar{y}) is

$$\frac{Y \text{ on } X}{Y - \bar{y} = \rho \frac{\sigma_y}{\sigma_x} (x - \bar{x})}$$

means
slope
coefficient of co-relation.

X on Y

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

Remark

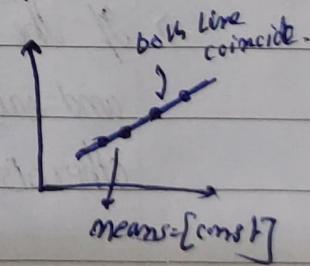
$$r = +1 \text{ or } -1$$

In a particular case of perfect co-relation that is +ve or -ve i.e., $r = +1$, then the line of regression for Y on X is :-

$$\frac{Y \text{ on } X}{Y - \bar{y} = \pm \frac{\sigma_y}{\sigma_x} (x - \bar{x})}$$

coincides
in case
of perfect
correlation
 $r = +1 \text{ or } -1$

$$\frac{X \text{ on } Y}{X - \bar{x} = \pm \frac{\sigma_x}{\sigma_y} (y - \bar{y})}$$



Since both lines of regression pass through the point (\bar{x}, \bar{y}) , so they can't be parallel, hence in this case of perfect correlation (+ve or -ve) the two lines will coincide.

* Regression coeff.

⇒ "b" → if it is the slope of the line of regression of Y on X is also called the coefficient of regression of Y on X .

$$b_{yx} = \gamma \frac{\sigma_y}{\sigma_x} = \text{Reg coeff of } Y \text{ on } X$$

$$b_{xy} = \gamma \frac{\sigma_x}{\sigma_y} = \text{Reg coeff of } X \text{ on } Y.$$

Q In a partially destroyed lab, record of an analysis of co-relation data. The following data are the only legible:-

$$\text{Var of } x = 9 \rightarrow \sigma_x^2 = 9, \sigma_x = 3.$$

$$\text{Regression eqn are } 8x - 10y + 66 = 0. \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$40x - 18y = 214$$

since the line of Reg passes through a point (\bar{x}, \bar{y}) so, means are:-

i) what are the mean values of x and y .

ii) Correlation coeff b/w x and y .

iii) σ_y .

$$\text{(i) } (\bar{x}, \bar{y}) = (13, 17)$$

$$8x - 10y = -66 \times 5$$

$$40x - 18y = 214$$

$$32y = 544$$

$$(ii) 8x - 10y + 66 = 0$$

$$40x - 18y = 214$$

$$y = \frac{544 + 36 - 17}{32} = 17$$

Y on X

X on Y

$$-10y = -66 - 8x$$

$$40x = 214 + 18y$$

$$\bar{y} = 17$$

$$y = \frac{66}{10} + \frac{8}{10}x$$

$$x = \frac{18}{40}y + \frac{214}{40}$$

$$8x - 10 \times 17 = -66$$

$$y = \frac{8x}{10} + \frac{66}{10}$$

$$b_{XY} = \gamma \frac{\sigma_x}{\sigma_y} = \frac{18}{40}$$

$$8x - 170 = -66$$

$$8x = 104$$

$$x = 13$$

$$\bar{x} = 13$$

$$b_{yx} = \gamma \frac{\sigma_y}{\sigma_x} = \frac{8}{10}$$

$$\gamma^2 = \frac{b_{yx} \cdot b_{xy}}{2 \cdot \frac{8}{10} \cdot \frac{18}{40}} = \frac{\gamma = \pm 0.6}{}$$

$$\text{iii) } \sigma_y = b_{yx} \cdot \frac{8}{10} = 0.6 \cdot \frac{\sigma_y}{\sigma_x}$$

7/09/08

(iii) σ_y $b_{yx} = 8/10$

$$b_{xy} = 18/40$$

$$\sigma_y =$$

$$r^2 \pm 0.6$$

$$b_{yx} = r \frac{\sigma_y}{\sigma_x}$$

$$\frac{8}{10} = \frac{0.8 \times \frac{8}{10}}{0.6 \sqrt{2}}$$

$$\boxed{\sigma_y = 4}$$

~~H.W~~ Q. Find the most likely price in Mumbai corresponding to the price of ₹70 at Kolkata from following :-

Mean $\rightarrow \bar{x}$ and \bar{y} 65

Kolkata (x) $\xrightarrow{\text{independent}}$ Mumbai (y)

67

$\xrightarrow{\text{dependent}} y = a + bx$

standard deviation $\sigma_x = 2.5$

$\sigma_y = 3.5$

Correlation Coeff (r) b/w the price of commodities in the two cities is 0.8.

$$\Rightarrow Y - \bar{Y} = \frac{\sigma_y}{\sigma_x} \frac{0.8}{\sigma_x} \frac{35}{65} (X - \bar{x})$$

$$Y - \bar{Y} = 0.8 \times \frac{7}{5} (X - 65)$$

$$Y - 67 = 0.8 \times \frac{7}{5} \times 8$$

$$Y = 5.6 + 67$$

$$\boxed{Y = 72.6}$$

* Properties of Correlation

1) The correlation coefficient is the geometric mean b/w the regression co-efficients.

$$\Leftrightarrow \boxed{r = \sqrt{b_{xy} \cdot b_{yx}}}$$

$a \& b$

gm $\rightarrow \sqrt{ab}$

geometric $\left\{ \begin{array}{l} am \rightarrow \frac{a+b}{2} \\ arithmetic \end{array} \right.$

$$b_{xy} \cdot b_{yx} = \left(r \frac{\sigma_y}{\sigma_x} \right) \cdot r \left(\frac{\sigma_x}{\sigma_y} \right)$$

$$= r^2$$

$$\sqrt{b_{xy} \cdot b_{yx}} = \sqrt{r^2} = r \text{ (proved)}$$

2) If one of the Regression co-efficient is greater than unity, other must be less than unity.

$$\Leftrightarrow b_{yx} > 1 \rightarrow \text{given}$$

$$\Rightarrow \frac{1}{b_{yx}} < 1 \quad \text{--- (1)}$$

$$r \text{ always lies } \rightarrow -1 \leq r \leq 1 \Rightarrow r^2 \leq 1$$

$$\Rightarrow b_{xy} \cdot b_{yx} \leq 1$$

$$b_{xy} \leq \frac{1}{b_{yx}}$$

from eq(1)

$$\boxed{b_{xy} \leq \frac{1}{b_{yx}} < 1}$$

$$\left. \begin{array}{l} b_{xy} \leq 1 \\ b_{yx} > 1 \end{array} \right\}$$

3) The modulus value of the arithmetic mean of regression coefficients is not less than the modulus value of correlation coefficient.

$$\Rightarrow \left| \frac{b_{xy} + b_{yx}}{2} \right| > |\gamma|$$

$$\Rightarrow \left| \frac{1}{2} \left(\gamma \frac{\sigma_x}{\sigma_y} + \gamma \frac{\sigma_y}{\sigma_x} \right) \right| > |\gamma|$$

$$\left| \frac{\sigma_y}{\sigma_x} + \frac{\sigma_x}{\sigma_y} \right| > 2$$

$$\left| \frac{\sigma_y^2 + \sigma_x^2}{\sigma_x \sigma_y} \right| > 2 \Rightarrow \sigma_y^2 + \sigma_x^2 > 2 \sigma_x \sigma_y$$

$$\Rightarrow \sigma_y^2 + \sigma_x^2 - 2 \sigma_x \sigma_y > 0$$

$$\Rightarrow (\sigma_x - \sigma_y)^2 > 0$$

which is true,
real quantity is always ≥ 0 .

4) Regression co-efficient are independent of change of origin but not of scale.

$$\hookrightarrow \text{Let } U = \frac{x-a}{h} \quad \& \quad V = \frac{y-b}{k}$$

$$\Rightarrow x = a + hU \quad \& \quad y = b + kV$$

where $a, b, h > 0$ & $k \neq 0$ are const.

$$b_{yx} = \gamma \frac{\sigma_y}{\sigma_x} = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} \cdot \frac{\sigma_y}{\sigma_x}$$

$$= \frac{\text{Cov}(X, Y)}{\sigma_x^2}$$

$$= \frac{hk \text{Cov}(U, V)}{\sigma_v^2}$$

∴ by

$$b_{XY} = \frac{h}{k} \frac{\text{Cov}(U, V)}{\sigma_v^2}$$

$$\boxed{b_{YX} = \frac{k}{h} \frac{\text{Cov}(U, V)}{\sigma_v^2}}$$

$$\text{or } \boxed{b_{YX} = \frac{k}{h} b_{UV}}$$

$$\boxed{b_{XY} = \frac{h}{k} b_{UV}}$$

5) Angle b/w two lines of regression is

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

* Case I \rightarrow If $\gamma=0$, $\tan \theta = \infty$

$$\Rightarrow \boxed{\theta = \frac{\pi}{2}}$$

Thus if two variables are un-correlated

i.e. $\boxed{\gamma=0}$, the line of regression perpendicular to each other.

* Case II \rightarrow $\gamma = \pm 1$, $\tan \theta = 0 \Rightarrow \boxed{\theta = 0 \text{ or } \pi}$

In this case, the two lines of regression either coincide or they are parallel to each other.

Remark \rightarrow Since both line of regression passes through the point (\bar{x}, \bar{y}) . so they can't be parallel. Hence, in the case of perfect correlation (+ve or -ve) the two lines of regression coincides.

* Proof of correlation coefficient property ②.

$$\text{Proof} \rightarrow \text{Let } U = \frac{X-a}{h} \quad \& \quad V = \frac{Y-b}{k}$$

$$\Rightarrow X = a + Uh \quad \& \quad Y = b + Kv \quad \text{--- ①}$$

Taking Expectation :-

$$\Rightarrow E[X] = a + E(U) \cdot h \quad \& \quad E[Y] = b + k \cdot E(V) \quad \text{--- ②}$$

① - ②

$$X - E[X] = h(U - E[U]) \quad \& \quad Y - E[Y] = k(V - E[V])$$

$$\begin{aligned} \text{Cov}(X, Y) &= \underbrace{(X - E[X])}_{\sim U_X} \underbrace{(Y - E[Y])}_{\sim V_Y} \\ &= h(U - E(U)) \cdot k(V - E(V)) \\ &= hk \cdot \{(U - E(U))(V - E(V))\} \end{aligned}$$

$$\boxed{\text{Cov}(X, Y) = hk \cdot \text{Cov}(U, V)} \quad \text{--- ①}$$

$$\begin{aligned} \sigma_x^2 &= E[(X - E(X))^2] \\ &= E[h(U - E(U))^2] \\ &= h^2 \cdot (U - E(U))^2 \\ &= h^2 \cdot \sigma_U^2 \end{aligned}$$

$$\begin{aligned} \sigma_y^2 &= E[(Y - E(Y))^2] \\ &= E[k(V - E(V))^2] \\ &= k^2 E[(V - E(V))^2] \\ &= k^2 \sigma_V^2 \end{aligned}$$

$$\boxed{\sigma_x^2 = h \cdot \sigma_U^2} \quad \text{--- ②}$$

$$\boxed{\sigma_y^2 = k \cdot \sigma_V^2} \quad \text{--- ③}$$

8/09/23

(UNIT-3) Some Discrete distributions

- The Bernoulli process.
- Binomial distribution and its Mgf [Moment Generating Funcⁿ] to find mean and Variance.
- Negative Binomial distribution and its Mgf.
- Geometric distribution and its Mgf.
- Poisson distribution & its Mgf.

$$\int_{m_0}^n p^r q^{n-r}$$

↑ prob of success
no. of trial ↓ prob of failure

* Bernoulli process

- An Experiment often consist of repeated trials, each with two possible outcome that may be denoted as success or failure.

every outcome can be seen as success or failure.

• Remark -

- 1) The Experiment consist of repeated trial. → Ex:- for tossing 2 coin.
 - 2) The probability of success is denoted by → "p",
which remains constant.
Heads can be
 $x = \{0, 1, 2\}$
↓
No. of successes.
- ⇒ The repeated trials are independent of each other.
- ⇒ The number of trials → "n" is always finite → for bernoulli;
infinite → for poisson.

* can be applied when 'n' is small.

* Binomial Distribution

A Bernoulli trial can result in success with probability "P" and a failure with probability "q", i.e., $[q = 1 - P]$, then the
 $[P + q = 1]$

probability distribution of the binomial RV 'x', the no. of success in 'n' independent trials is

$$b(x; n, p) = {}^n C_x p^n \cdot q^{n-x}$$

Q Three items are selected at random from a manufacturing process, inspected and classified as defective or non-defective. Find the probability distribution for the number of defective items, the assuming that 25% of items are defective.

⇒ Let X :- defective. $n = 3$

$$\frac{25\% \text{ of } 3}{100} = \frac{25 \times 3}{100} = \frac{75}{100}$$

$$X = \{0, 1, 2, 3\}$$

$$p = \frac{25}{100} = \frac{1}{4}$$

$$b(x; n, p) = {}^3 C_x \left(\frac{1}{4}\right)^x \cdot \left(\frac{3}{4}\right)^{n-x}$$
$$x = \{0, 1, 2, 3\}$$
$$x=0 = {}^3 C_0 \cdot \left(\frac{1}{4}\right)^0 \cdot \left(\frac{3}{4}\right)^3$$
$$= 1 \cdot \frac{1}{4} \cdot \frac{27}{256}$$

Probability of defective.

$$x=1 = {}^3 C_1 \cdot \left(\frac{1}{4}\right)^1 \cdot \left(\frac{3}{4}\right)^2$$

$$x=2 = {}^3 C_2 \cdot \left(\frac{1}{4}\right)^2 \cdot \left(\frac{3}{4}\right)^1$$

$$x=3 = {}^3 C_3 \cdot \left(\frac{1}{4}\right)^3 \cdot \left(\frac{3}{4}\right)^0$$

$$b = 1 \cdot \left(\frac{1}{4}\right)^0 \cdot \frac{27}{256} + 3 \cdot \left(\frac{1}{4}\right)^1 \cdot \left(\frac{3}{4}\right)^2 + 3 \cdot \left(\frac{1}{4}\right)^2 \cdot \left(\frac{3}{4}\right)^1 + 1 \cdot \left(\frac{1}{4}\right)^3$$
$$= \frac{27}{256} + \frac{27}{64} + \frac{9}{64} + \frac{1}{64} = \frac{27+108+36+1}{256} = \frac{172}{256} = \underline{\underline{0.675}}$$

11/09/23

Q The prob that a patient recovers from a rare blood disease is 0.4. If 15 people are known to have contracted this disease. What is the prob atleast of:

i) atleast 10 survive.

(ii) Exactly 5 survive.

(iii) From 3 to 8 patient are survive.

$$n=15$$

Let $X \sim$ Survive.

$$x = \{0, \dots, 15\}$$

$$\begin{array}{|l} \text{prob} = (0.4)^x \cdot p \\ \text{of recovery.} \end{array}$$

$$q = 0.6$$

$$b(x; n, p) = \binom{n}{x} \cdot p^x \cdot q^{n-x} = {}^{15}C_x \cdot p^x \cdot q^{15-x}$$

(i) $P(X \geq 10) \Rightarrow x = 10, 11, 12, 13, 14, 15.$

$$b(x; 15, 0.4) = {}^{15}C_x \cdot p^x \cdot q^{15-x}$$

$$b_{10} \Rightarrow x=10 = {}^{15}C_{10} \cdot p^{10} \cdot q^5$$

$$\frac{15!}{10!5!} \cdot \left(\frac{4}{10}\right)^{10} \cdot \left(\frac{6}{10}\right)^5$$

$$x=11 = {}^{15}C_{11} \cdot p^{11} \cdot q^4$$

$$\frac{15!}{14!1!} \cdot \frac{1}{10!5!} \cdot \frac{1}{4!3!2!1!}$$

$$x=12 = {}^{15}C_{12} \cdot p^{12} \cdot q^3$$

$$\frac{15!}{13!2!} \cdot \frac{1}{10!5!} \cdot \frac{1}{3!2!1!}$$

$$555$$

$$x=13 \rightarrow {}^{15}C_{13} \cdot p^{13} \cdot q^2$$

$$\frac{15!}{12!3!} \cdot \frac{1}{10!5!} \cdot \frac{1}{2!1!}$$

$$x=14 \rightarrow {}^{15}C_{14} \cdot p^{14} \cdot q^1$$

$$\frac{15!}{13!2!} \cdot \frac{1}{10!5!} \cdot \frac{1}{1!}$$

$$x=15 \rightarrow {}^{15}C_{15} \cdot p^{15} \cdot q^0$$

②

* Mean of binomial

$$\hookrightarrow \boxed{\mu = np}$$

* Variance of binomial

$$\hookrightarrow \boxed{\text{Var} = npq}$$

* (Moment Generating function) = MGF.

↳ The MGF of a RV "X" about origin having the probability function $f(x)$ is given by

$$M_X(t) = E[e^{tx}] = \begin{cases} \int e^{tx} f(x) dx, & \text{continuous} \\ \sum e^{tx} f(x), & \text{discrete.} \end{cases}$$

$$E[x] = \sum x f(x)$$

$$M_X(t) = E[e^{tx}] = E\left[1 + \frac{e^{tx}}{1!} + \frac{t^2 e^{tx}}{2!} + \dots + \frac{t^r e^{tx}}{r!} + \dots\right]$$

$$= E\left[1 + tx + \frac{t^2 x^2}{2!} + \frac{t^3 x^3}{3!} + \dots + \frac{t^r x^r}{r!}\right] \quad e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots$$

$$= 1 + t E[X] + \frac{t^2 E[X^2]}{2!} + \dots + \frac{t^r E[X^r]}{r!} + \dots$$

moment of r^{th} term

$$= 1 + t M'_1 + \frac{t^2 M'_2}{2!} + \frac{t^3 M'_3}{3!} + \dots + \frac{t^r M'_r}{r!} + \dots$$

$$= \sum_{r=0}^{\infty} \frac{t^r}{r!} \cdot M'_r$$

$$\Rightarrow \boxed{M_X(t) = \sum_{r=0}^{\infty} \frac{t^r}{r!} M'_r}$$

or

$$\boxed{M_X^{(r)} = \left[\frac{d^r}{dt^r} M_X(t) \right]_{t=0}}$$

Find mean and Variance of a binomial distribution by using MGF :-

Let 'X' be a binomial R.V. such that,

$$b(n; n, p) = {}^n C_x \cdot p^n \cdot q^{n-x} ; x = 0, 1, 2, \dots, n$$

$$M_X(t) = E[e^{tx}] = \sum e^{tx} \cdot {}^n C_x \cdot p^x \cdot q^{n-x}$$

$$= \sum {}^n C_x \cdot (pe^t)^x \cdot q^{n-x}$$

≡ (1)

$$(1+x)^n = 1 + nx \frac{n}{1!} + n \frac{(n-1)}{2!} x^2 + n \frac{(n-1)(n-2)}{3!} x^3 + \dots$$

$$(n+a)^n = {}^n C_0 a^n x^0 + {}^n C_1 a^{n-1} x^1 + \dots$$

$$= \sum_{x=0}^n (pe^t + q)^n$$

Mean ⇒ $\mu = E[X] = \mu'$ $p+q=1$

$$\mu' = \left| \frac{d}{dt} M_X(t) \right|_{t=0} = \left| \frac{d}{dt} (q + pe^t)^n \right|_{t=0}$$

$$\mu' = \left| \frac{d}{dt} (q + pe^t)^n \right|_{t=0} = \left| n(q+p)e^{nt} \right|_{t=0}$$

$$= \left| n(q+pe^t)^{n-1} \cdot pe^t \right|_{t=0} = \left| n(1)^{n-1} \cdot p \right|$$

$E[X] = np$

$$M_2' = E[X^2]$$

$$\downarrow \frac{d^r}{dt^r} M_X(t)$$

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

↓ mean

$$E[X^2] = M_2'$$

$$= \left. \frac{d^2}{dt^2} M_X(t) \right|$$

$$= \left. \frac{d^2}{dt^2} (a + be^t)^n \right|_{t=0}$$