

Probability and Statistics

- ① Discriptive Statistics
- ② Correlation and Regression
- ③ Probability Distributions
- ④ Sampling Theory
- ⑤ Tests of Hypothesis

Data Science collection of data
Statistics primary data
Population secondary data
Sample

C.I	f	c.f
20-30	1	1
30-40	3	4
40-50	2	6
50-60	3	9
60-70	5	14
70-80	4	18
80-90	2	20
90-100	1	21

Line diagram

Bar diagram

Rectangular diagram

Histogram

Pie diagram

Measures of Central Tendency

- (i) Mean / Arithmetic Mean / Agrigrade
- (ii) Median
- (iii) Mode
- (iv) Geometric Mean
- (v) Harmonic Mean

① Arithmetic Mean

(i) Individual Series : If $x_1, x_2, x_3, \dots, x_n$ are n observations, then

$$A.M = \frac{1}{n} \sum_{i=1}^n x_i$$

Eg: 44 65 70 68 64 85

$$A.M = \frac{44+65+70+68+64+85}{6} = 66$$

② Discrete Frequency Distribution

If $x_1, x_2, x_3, \dots, x_n$ are n observations with corresponding frequencies

$f_1, f_2, f_3, \dots, f_n$

x_i	f_i	$x_i f_i$
x_1	f_1	$f_1 x_1$
x_2	f_2	$f_2 x_2$
x_3	f_3	$f_3 x_3$
\vdots	\vdots	\vdots
x_n	f_n	$f_n x_n$
	N	$\sum f_i x_i$

$$A.M = \frac{1}{N} \sum f_i x_i$$

where $n = \sum f_i$

Eg:

x_i	f_i	$x_i f_i$
1	5	5
2	9	18
3	12	36
4	17	68
5	14	70
6	10	60
7	6	42
	73	299

$$A.M = \frac{1}{N} \sum f_i x_i$$

$$= \frac{1}{73} \times 299 = 4.09$$

C.I	f_i	Middle points x_i	$f_i x_i$
$I_0 - I_1$	f_1	x_1	$f_1 x_1$
$I_1 - I_2$	f_2	x_2	$f_2 x_2$
$I_2 - I_3$	f_3	x_3	$f_3 x_3$
\vdots	\vdots	\vdots	\vdots
$I_{n-1} - I_n$	f_n	x_n	$f_n x_n$
	$\sum f_i$		$\sum f_i x_i$

$$A \cdot M = \frac{\sum f_i x_i}{\sum f_i}$$

→ Find the mean of the following frequency distribution

Marks	0-10	10-20	20-30	30-40	40-50	50-60
No. of Students	12	18	27	20	17	6

C.I	f_i	x_i	$f_i x_i$
0-10	12	5	60
10-20	18	15	270
20-30	27	25	675
30-40	20	35	700
40-50	17	45	765
50-60	6	55	330

$$\sum f_i = 100$$

$$\sum f_i x_i = 2800$$

$$A \cdot M = \frac{\sum f_i x_i}{\sum f_i} = \frac{2800}{100} = 28$$

Deviation Method

case-(i): $d_i = x_i - A$ where 'A' is Assumed Mean

$$A \cdot M = A + \frac{\sum f_i d_i}{\sum f_i}$$

case-(ii): $d_i = \frac{x_i - A}{h}$

$$A \cdot M = A + \frac{\sum f_i d_i}{N} \times h$$

where $N = \sum f_i$

$h = \text{length of the class}$

→ Find the A.M of the following distribution by using Deviation Method.

Marks	0-10	10-20	20-30	30-40	40-50	50-60
No. of Students	12	18	27	20	17	6

C.I	f _i	x _i	d _i = $\frac{x_i - A}{h}$	f _i d _i
0-10	12	5	-3	-36
10-20	18	15	-2	-36
20-30	27	25	-1	-27
30-40	20	35	0	0
40-50	17	45	1	17
50-60	6	55	2	12

$$\sum f_i = N = 100$$

$$\sum f_i d_i = -70$$

$$A.M = A + \frac{\sum f_i d_i}{N} \times h$$

$$A.M = 35 + \frac{-70}{100} \times 10$$

$$\boxed{A.M = 28}$$

Median

(i) Individual Series

'n' no. of observations

If n is odd,

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ term}$$

If n is even,

$$\text{Median} = \text{Avg of } \left(\frac{n}{2} \right)^{\text{th}} \text{ term and } \left(\frac{n+1}{2} \right)^{\text{th}} \text{ term}$$

① Find the median of the following data

4 7 12 15 13 5 4

Ascending Order : 4 4 5 7 12 13 15

$n=7$ (no. of observations)

$$\text{Median} = \left(\frac{n+1}{2}\right)^{\text{th}} \text{ term}$$

$$= \frac{7+1}{2}^{\text{th}} \text{ term}$$

$$= 4^{\text{th}} \text{ term}$$

$$= 15$$

② Find the median of the following data

4 6 12 15 20 24 2 10 35 40

Ascending order: 2 4 6 10 12 15 20 24 35 40

$n=10$ (no. of observations)

$$\text{Median} = \frac{\left(\frac{n}{2}\right)^{\text{th}} + \left(\frac{n}{2}+1\right)^{\text{th}}}{2}$$

$$= \frac{5^{\text{th}} + 6^{\text{th}}}{2}$$

$$= \frac{12+15}{2}$$

$$= \frac{27}{2}$$

$$= 13.5$$

(ii) Discrete Frequency Distribution

x	0	1	2	3	4	5	6	7	8	9
f	3	5	12	45	61	22	29	51	23	3

x	f	cdf	
0	3	3	
1	5	8	$N = 254$
2	12	20	
3	45	65	$\frac{N}{2} = \frac{254}{2}$
4	61	126	$= 127$
5	22	148	
6	29	177	Median = 5
7	51	228	
8	23	251	
9	3	254	

(iii) Continuous Frequency Distribution

$$\text{Median} = l + \frac{\frac{N}{2} - m}{f} \times h$$

where l = lower limit of the median class

f = frequency of the median class

$$N = \sum f$$

m = cumulative frequency of the class preceding the median class

h = length of the class interval

→ Find the median of the following frequency distribution.

C.I	0-10	10-20	20-30	30-40	40-50	50-60
f	12	18	27	20	17	6

C.I	f	cf	
0-10	12	12	Here, $l = 20$
10-20	18	30	$f = 27$
20-30	27	57	$m = 30$
30-40	20	77	$= 10$
40-50	17	94	
50-60	6	100	$\text{Median} = l + \frac{\frac{N}{2} - m}{f} \times h$

$$N = 100$$

$$= 20 + \frac{50-30}{27} \times 10$$

$$\frac{N}{2} = 50$$

$$= 20 + \frac{20}{27} \times 10$$

$$= 27.407$$

Mode: Most frequently occurring value in a given series

(i) Individual Series

3, 4, 3, 5, 3, 7, 7, 7, 9, 5, 3, 10, 3

Mode = 3

(ii) Discrete Frequency Distribution

x	1	2	3	4	5	6
f	5	15	25	15	20	6

Mode = 3

(iii) Continuous Frequency Distribution

$$\text{Mode} = l + \frac{f_1 - f_0}{(f_1 - f_0) - (f_2 - f_1)} \times h$$

$$\text{Mode} = l + \frac{f_1 - f_0}{af_1 - f_0 - f_2} \times h$$

where l = lower limit of the modal class

f_1 = frequency of the modal class

f_0 = frequency of the class preceding the modal class

f_2 = frequency of the class succeeding the modal class

h = length of the class interval

→ Find the mode of the following frequency distribution.

C.I	0-10	10-20	20-30	30-40	40-50	50-60	60-70
f	5	6	8	12	15	5	3

Here $l = 40$

$f_1 = 15$

$f_0 = 12$

$f_2 = 5$

$h = 10$

$$\text{Mode} = l + \frac{f_1 - f_0}{(f_1 - f_0) - (f_2 - f_1)} \times h$$

$$= 40 + \frac{3}{3+10} \times 10$$

$$= 40 + \frac{30}{13}$$

$$= 42.307$$

Geometric Mean

(i) Individual Series

If $x_1, x_2, x_3, \dots, x_n$ are 'n' observations, then

$$G.M = \sqrt[n]{x_1 \cdot x_2 \cdot x_3 \cdots x_n}$$

$$\log G.M = \log(x_1 \cdot x_2 \cdots x_n)^{1/n}$$

$$= \frac{1}{n} \log(x_1 \cdot x_2 \cdots x_n)$$

$$= (\log x_1 +$$

$$= \frac{1}{n} \sum_{i=1}^n \log x_i$$

$$G.M = n \text{ Antilog} \left(\frac{1}{n} \sum_{i=1}^n \log x_i \right)$$

For frequency distribution,

$$G.M = \text{Antilog} \left(\frac{1}{N} \sum_{i=1}^N f_i \log x_i \right)$$



Harmonic Mean: It is defined as "the reciprocal of avg of reciprocals of the given set of observations"

(i) Individual Series:

If $x_1, x_2, x_3, \dots, x_n$ are 'n' observations, then their reciprocals are

$$\frac{1}{x_1}, \frac{1}{x_2}, \dots, \frac{1}{x_n}$$

$$A.M = \frac{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}{n} = \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}$$

$$H.M = \frac{1}{\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}}$$

$$= \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

For a frequency distribution,

$$H.M = \frac{N}{\sum_{i=1}^n f_i \left(\frac{1}{x_i} \right)} \quad (\text{where } N = \sum_{i=1}^n f_i)$$

→ Find the H.M. of the set of observations: 2, 4, 8, 12, 16, 24

x_i	$\frac{1}{x_i}$	$n = 6$
2	0.5	$H.M = \frac{n}{\sum_{i=1}^n \left(\frac{1}{x_i}\right)}$
4	0.25	
8	0.125	
12	0.083	$= \frac{6}{1.062}$
16	0.062	$= 5.649$
24	0.042	
		<u>1.062</u>

→ Find G.M and H.M of the following frequency distribution.

C.I	0-10	10-20	20-30	30-40	40-50	50-60	60-70
f	4	7	13	21	15	25	20

C.I	f_i	x_i	$\log x_i$	$f_i \log x_i$	$\frac{1}{x_i}$	$f_i \left(\frac{1}{x_i}\right)$
0-10	4	5	0.6989	2.7956	0.200	0.800
10-20	7	15	1.1760	8.2320	0.066	0.462
20-30	13	25	1.3979	18.17	0.040	0.520
30-40	21	35	1.5440	32.4240	0.028	0.588
40-50	15	45	1.6532	24.7980	0.022	0.330
50-60	25	55	1.7403	43.5075	0.018	0.450
60-70	20	65	1.8129	36.2580	0.015	0.300
	<u>105</u>		<u>10.0232</u>	<u>166.1878</u>		<u>3.450</u>

$$\text{Here } N = \sum f_i = 105$$

$$\begin{aligned}
 G.M &= \text{Antilog} \left(\frac{1}{N} \sum_{i=1}^n f_i \log x_i \right) & H.M &= \frac{N}{\sum_{i=1}^n f_i \left(\frac{1}{x_i} \right)} \\
 &= \text{Antilog} \left(\frac{1}{105} \times 166.1878 \right) & &= \frac{105}{3.450} \\
 &= \text{Antilog} (1.5327) & &= 30.4347 \\
 &= 38.2596
 \end{aligned}$$

Measures of Dispersion:

- (i) Range
- (ii) Quartile Deviation
- (iii) Mean Deviation
- (iv) Standard Deviation

Range:

$$\text{Range} = \text{Max value} - \text{Min value}$$

Let $X \rightarrow x_1, x_2, x_3, \dots, x_n$

$$\text{Range} = X_{\max} - X_{\min}$$

$$\text{Coeff of Range} = \frac{X_{\max} - X_{\min}}{X_{\max} + X_{\min}}$$

→ Find the range of the given set of observations : 48, 82, 64, 72, 66, 68

$$\begin{aligned}\text{Range} &= X_{\max} - X_{\min} \\ &= 82 - 48 \\ &= 34\end{aligned}$$

$$\begin{aligned}\text{Coeff of Range} &= \frac{82 - 48}{82 + 48} = \frac{34}{130} \\ &= 0.2615\end{aligned}$$

Quartile Deviation (Q.D)

$$Q.D = \frac{Q_3 - Q_1}{2}$$

~~Q₁~~ → First Quartile

Q₃ → Third Quartile

$$\text{where, } Q_i = l + \frac{\frac{N}{4} - m}{f} \times c$$

l = lower limit of the first quartile class

f = frequency of the first quartile class

m = cf of the class preceding the first quartile class

c = length of the class interval.

$$Q_3 = l + \frac{\frac{3N}{4} - m}{f} \times c$$

l = lower limit of the third quartile class

f = frequency of the third quartile class

m = cf of the class preceding the third quartile class

c = length of the class interval

→ Find the quartile deviation of the following frequency distribution.

C.I	0-15	15-30	30-45	45-60	60-75	75-90	90-105
f	8	26	30	45	20	17	4

C.I	f	cf
0-15	8	8
15-30	26	34
30-45	30	64
45-60	45	109
60-75	20	129
75-90	17	146
90-105	4	150

$$N = 150$$

$$\frac{N}{4} = \frac{150}{4} = 37.5$$

$$l_1 = 30, f_1 = 30, m_1 = 34, c = 15$$

$$Q_1 = l_1 + \frac{\frac{N}{4} - m_1}{f_1} \times c$$

$$= 30 + \frac{37.5 - 34}{30} \times 15$$

$$= 31.75$$

$$N = 150$$

$$\frac{3N}{4} = \frac{450}{4} = 112.5$$

$$l_2 = 60, f_2 = 20, m_2 = 109, c = 15$$

$$Q_3 = l_2 + \frac{\frac{3N}{4} - m_2}{f_2} \times c$$

$$= 60 + \frac{112.5 - 109}{20} \times 15$$

$$= 62.625$$

$$Q.D = \frac{Q_3 - Q_1}{2}$$

$$= \frac{62.625 - 31.75}{2}$$

$$= 15.4375$$

Mean Deviation

(i) Individual Series

If $x_1, x_2, x_3, \dots, x_n$ are 'n' observations, then

$$M.D = \frac{1}{n} \sum_{i=1}^n |(x_i - A)|$$

where 'A' is any aggregate

For frequency distribution,

$$M.D = \frac{1}{n} \sum_{i=1}^n f_i |x_i - A|$$

→ Find the mean deviation ^{about} from the mean to the following data.

Marks less than	80	70	60	50	40	30	20	10
No. of students	100	80	70	60	32	20	13	5

C.I	f_i	x_i	$f_i x_i$	$ x_i - \bar{x} $	$f_i x_i - \bar{x} $
0-10	5	5	25	42	$\sum f_i = 100$
10-20	8	15	120	32	$\sum f_i x_i = 4700$
20-30	7	25	175	22	
30-40	12	35	420	12	$\text{Mean} = \frac{\sum f_i x_i}{N}$
40-50	28	45	1260	2	
50-60	10	55	550	8	
60-70	10	65	650	18	
70-80	20	75	1500	28	
	<u>100</u>		<u>4700</u>		

$$\sum f_i = 100, \sum f_i x_i = 4700$$

$$\text{Mean} = \frac{1}{N} \sum f_i x_i = \frac{4700}{100} = 47$$

$$\sum f_i |(x_i - \bar{x})| = 1640$$

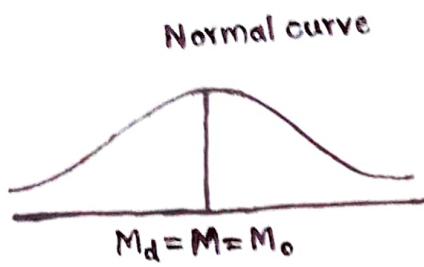
$$M.D = \frac{1}{n} \sum_{i=1}^n f_i |x_i - \bar{x}|$$

$$= \frac{1640}{100}$$

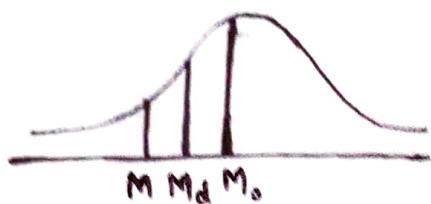
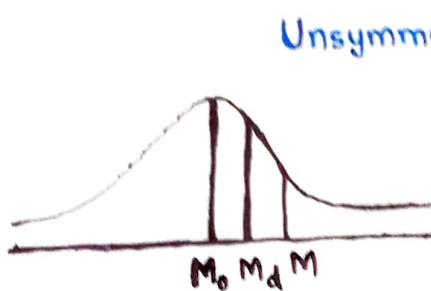
$$= 16.4$$

Skewness

Lack of Symmetry



Symmetric curve



Measures

$$S_K = \text{Mean} - \text{Median}$$

$$= \text{Mean} - \text{Mode}$$

$$S_K = (Q_3 - M_d) - (M_d - Q_1)$$

$$= Q_3 + Q_1 - 2M_d$$

$$\text{Mode} = 3 \text{Median} - 2 \text{Mean}$$

$$M_o = 3M_d - 2M$$

Absolute value of skewness

Karl Pearson's

↓

$$\text{Coeff of skewness} = \frac{M - M_o}{\text{S.D.}}$$

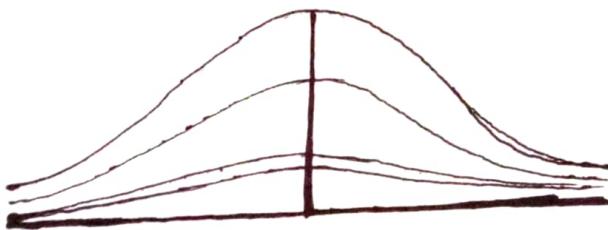
$$= \frac{M - (3M_d - 2M)}{\sigma} = \frac{3(M - M_d)}{\sigma}$$

Bowley's coeff of skewness

$$= \frac{(Q_3 - M_d) - (M_d - Q_1)}{(Q_3 - M_d) + (M_d - Q_1)}$$

Kurtosis

Convexity of the curve



Measures the flatness of the frequency distribution

A → Leptokurtic (sharp)

B → Mesokurtic (Normal)

C → platykurtic (flat)

Coeff of kurtosis whose measures are represented by β_2 and γ_2

$$\beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{\mu_4}{\sigma^4}$$

$$\gamma_2 = \beta_2 - 3$$

where $\mu_r = r^{\text{th}} \text{ moment about mean}$

$$\mu_r = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^r$$

Correlation and Regression

Correlation

Positive Correlation

Negative Correlation

No Correlation

Simple Correlation

Multiple Correlation

Partial Correlation

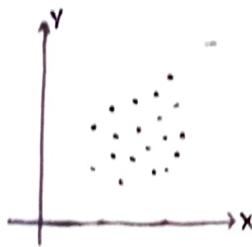
Linear Correlation

$$x: \underline{3} \quad \underline{5} \quad \underline{7} \quad \underline{9} \quad \underline{11}$$

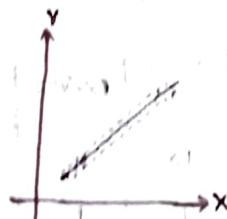
Curvilinear Correlation

$$y: \underline{7} \quad \underline{11} \quad \underline{15} \quad \underline{19} \quad \underline{23}$$

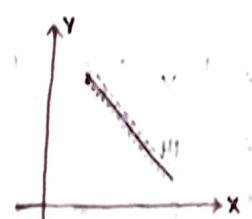
Scatter diagram



$$\bar{x} = \frac{a}{4} = \frac{2}{4} = \frac{3}{4} = \frac{a}{4}$$



Positive Correlation



Negative Correlation

Karl Pearson's Coeff of Correlation

$r(x, y)$, r_{xy}

co-variance

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

$$r = \frac{\frac{1}{n} \sum_{i=1}^n (x - \bar{x})(y - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2}}$$

$$r = \frac{\frac{1}{n} \sum_{i=1}^n (x - \bar{x})(y - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \bar{y})^2}}$$

Properties

$$\gamma = \gamma(x, y)$$

(I) $-1 \leq \gamma \leq 1$

If $0 < \gamma < 1$, then x, y are positively correlated

If $-1 < \gamma < 0$, then x, y are negatively correlated

If $\gamma = 0$, then there is no correlation

If $\gamma = 1$, then \exists perfect positive correlation b/w x and y .

If $\gamma = -1$, then \exists perfect negative correlation b/w x and y .

(II) Correlation coeff is independent of change of scale and origin

→ Find the correlation coeff b/w two variables X and Y follows

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

X	Y	$\frac{X-\bar{X}}{=X-10}$	$\frac{Y-\bar{Y}}{=Y-9}$	$(X-\bar{X})(Y-\bar{Y})$	$(X-\bar{X})^2$	$(Y-\bar{Y})^2$
12	14	2	5	10	4	25
9	8	-1	-1	1	1	1
8	6	-2	-3	6	4	9
10	9	0	0	0	0	0
11	11	1	2	2	1	4
13	12	3	3	9	9	9
7	3	-3	-6	18	9	36
<u>70</u>	<u>63</u>			<u>46</u>	<u>88</u>	<u>84</u>

$$\Sigma X = 70, \Sigma Y = 63, n = 7$$

Coeff of correlation

$$\bar{X} = \frac{1}{n} \sum X = \frac{70}{7} = 10$$

$$\gamma = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$$\bar{Y} = \frac{1}{n} \sum Y = \frac{63}{7} = 9$$

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

$$\sum (X-\bar{X})(Y-\bar{Y}) = 46$$

$$\sigma_x = \sqrt{\frac{88}{6}}$$

$$\sum (X-\bar{X})^2 = 88$$

$$\gamma = \frac{46}{\sqrt{88} \sqrt{84}} = 0.948 \approx 0.95$$

$$\sum (Y-\bar{Y})^2 = 84$$

∴ The two variables x and y are positively correlated.

→ Find if there is any significant correlation b/w the heights & weights given below.

Heights (in inches)	57	59	62	63	64	65	55	58	57
Weights (in pounds)	113	117	126	126	130	129	111	116	112

X	Y	(X - \bar{X})	(Y - \bar{Y})	(X - \bar{X})(Y - \bar{Y})	$(X - \bar{X})^2$	$(Y - \bar{Y})^2$
57	113	-3	-7	21	9	49
59	117	-1	-3	3	1	9
62	126	2	6	12	4	36
63	126	3	6	18	9	36
64	130	4	10	40	16	100
65	129	5	9	45	25	81
55	111	-5	-9	45	25	81
58	116	-2	-4	8	4	16
57	112	-3	-8	24	9	64
<u>540</u>	<u>1080</u>			<u>216</u>	<u>102</u>	<u>472</u>

$$\sum X = 540, \sum Y = 1080, n = 9$$

Coeff of correlation

$$\bar{X} = \frac{1}{n} \sum X = \frac{540}{9} = 60$$

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

$$\bar{Y} = \frac{1}{n} \sum Y = \frac{1080}{9} = 120$$

$$r = \frac{216}{\sqrt{102} \sqrt{472}}$$

$$\sum (X - \bar{X})(Y - \bar{Y}) = 216$$

$$r = 0.95$$

$$\sum (Y - \bar{Y})^2 = 472$$

→ When deviations are taken from Assumed Mean

$$r = \frac{\sum (X - A)(Y - B)}{\sqrt{\sum (X - A)^2} \sqrt{\sum (Y - B)^2}}$$

→ Find the correlation coeff for the following paired data.

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

X	Y	$\bar{x} = X - 35$	$\bar{y} = Y - 31$	xy	x^2	y^2
28	23	-7	-8	56	49	64
41	34	6	3	18	36	9
40	32	5	2	10	25	4
38	31	3	3	9	9	9
35	30	0	-1	0	0	1
33	26	-2	-5	10	4	25
40	28	5	-3	15	25	9
32	31	-3	0	-6	9	0
36	36	1	5	5	1	25
33	38	-2	7	-14	4	64
<u>356</u>	<u>313</u>	<u>6</u>	<u>3</u>	<u>79</u>	<u>162</u>	<u>195</u>

$$\sum X = 356, \sum Y = 313, n = 10$$

$$\bar{x} = \frac{1}{n} \sum X = \frac{356}{10} = 35.6$$

$$\bar{y} = \frac{1}{n} \sum Y = \frac{313}{10} = 31.3$$

$$\sum x = 6, \sum y = 3, \sum xy =$$

$$\sum x^2 = 162, \sum y^2 = 195$$

Coeff of correlation

$$r = \frac{\sum xy - \frac{1}{n} \sum x \sum y}{\sqrt{\sum x^2 - \frac{1}{n} (\sum x)^2} \sqrt{\sum y^2 - \frac{1}{n} (\sum y)^2}}$$

$$r = \frac{79 - \frac{1}{10}(6)(3)}{\sqrt{162 - \frac{1}{10}(356)^2} \sqrt{195 - \frac{1}{10}(313)^2}}$$

$$r = \frac{79 - 1.8}{\sqrt{162 - 3.6} \sqrt{195 - 0.9}}$$

$$r = \frac{77.2}{\sqrt{158.4} \sqrt{194.81}}$$

$$r = \frac{77.2}{(12.58)(13.93)}$$

$$r = \frac{77.2}{175.23}$$

$$r = 0.44$$

Spearman's Rank Correlation Coefficient

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

where $d_i = x_i - y_i$

x_i denotes ranks of X

y_i denotes ranks of Y

$$\text{For repeated ranks, } \rho = 1 - \frac{6 \sum d_i^2 + C.F.}{n(n^2-1)}$$

where $C.F. = \text{correction factor}$

$$C.F. = \frac{m(m^2-1)}{12} \quad (m = \text{no. of times the obs is repeated})$$

If any obs of a given series is repeated m times, then

$$C.F. = \sum_{i=1}^n \frac{m_i(m_i^2-1)}{12}$$

→ random sample of five college students is selected & their marks in Mathematics & Statistics are found to be as follows;

Math	85	60	73	40	90
Stat	3	75	65	50	80

x	y	x_i	y_i	$d_i = x_i - y_i$	d_i^2
85	93	2	1	1	1
60	75	4	3	1	1
73	65	3	4	-1	1
40	50	5	5	0	0
90	80	1	2	-1	1
					<u>4</u>

Here $n=5$

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} = 1 - \frac{6 \times 4}{5(25)} = 1 - \frac{1}{5} = 1 - 0.2 = 0.8$$

→ Find the rank correlation coefficient for the following data.

X	68	64	75	50	64	80	75	4	55	64
y	62	58	68	45	81	60	68	48	50	70

X	y	x_i	y_i	$d_i = x_i - y_i$	d_i^2
68	62	4	5	-1	1
64	58	6	7	-1	1
75	68	2.5	3.5	-1	1
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	60	8	8	0	0
64	70	6	4	2	16
					72
					33

Here $n=10$

In X-Series : 75 repeated 2 times

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

In Y-Series : 68 repeated 2 times

$$\rho = \frac{6(72+3)}{10 \times 99}$$

$$C.F = \frac{2(2^2-1)}{12} + \frac{3(3^2-1)}{12} + \frac{2(2^2-1)}{12}$$

$$\rho = 1 - 0.4545$$

$$\rho = 0.545$$

$$C.F = \frac{1}{2} + 2 + \frac{1}{2}$$

$$C.F = 3$$

Regression:

Linear Regression

Regression line of Y on X

$$(Y - \bar{Y}) = b_{yx}(X - \bar{X})$$

$b_{yx} \rightarrow$ regression coeff of
 Y on X

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

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

$$b_{yx} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{\sum (X - \bar{X})^2}$$

Regression line of X on Y

$$(X - \bar{X}) = b_{xy}(Y - \bar{Y})$$

$b_{xy} \rightarrow$ regression coeff of
 X on Y

$$b_{xy} = r \cdot \frac{\sigma_x}{\sigma_y}$$

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

$$b_{xy} = \frac{\sum (Y - \bar{Y})(X - \bar{X})}{\sum (Y - \bar{Y})^2}$$

→ Find the equations of the regression lines.

X	46	42	44	40	43	41	45
Y	40	38	36	35	39	37	41

X	Y	$X - \bar{X}$ $= X - 43$	$Y - \bar{Y}$ $= Y - 38$	$(X - \bar{X})(Y - \bar{Y})$	$(X - \bar{X})^2$	$(Y - \bar{Y})^2$
46	40	3	2	6	9	4
42	38	-1	0	0	1	0
44	36	1	-2	-2	1	4
40	35	-3	-3	9	9	9
43	39	0	1	0	0	1
41	37	-2	-1	2	4	1
45	41	2	3	6	4	9
301	266			21	28	28

$$\sum X = 301, \sum Y = 266, n = 7$$

$$\sum (X - \bar{X})(Y - \bar{Y}) = 21$$

$$\bar{X} = \frac{\sum X}{n} = \frac{301}{7} = 43$$

$$\sum (X - \bar{X})^2 = 28$$

$$\bar{Y} = \frac{\sum Y}{n} = \frac{266}{7} = 38$$

$$\sum (Y - \bar{Y})^2 = 28$$

Eqn to the regression line of y on x is

$$(y - \bar{y}) = b_{yx}(x - \bar{x})$$

$$b_{yx} = r \cdot \frac{\sigma_y}{\sigma_x} = \frac{\sum(x - \bar{x})(y - \bar{y})}{\sum(x - \bar{x})^2} = \frac{21}{28} = \frac{3}{4} = 0.75$$

Regression line of y on x is

$$(y - 38) = (0.75)(x - 43)$$

$$y = (0.75)(38 - 32.5)$$

$$y = (0.75) \cancel{(38 - 32.5)}$$

→ Find the regression equations b/w the indices of cotton & wool which are given below for the 12 months of a year.

Cotton(x)	78	77	85	88	87	82	81	77	76	83	97	93
Wool(y)	84	82	82	85	89	90	88	92	83	89	98	99

$$\sum x = 1004, \sum y = 1061, n = 12$$

$$\sum(x - \bar{x})(y - \bar{y}) = 288.66$$

$$\bar{x} = \frac{1}{n} \sum x = \frac{1004}{12} = 83.7$$

$$\sum(x - \bar{x})^2 = 486.68$$

$$\bar{y} = \frac{1}{n} \sum y = \frac{1061}{12} = 88.4$$

$$\sum(y - \bar{y})^2 = 362.92$$

X	Y	$X - \bar{X}$ $= X - 83.7$	$Y - \bar{Y}$ $= Y - 88.4$	$(X - \bar{X})(Y - \bar{Y})$	$(X - \bar{X})^2$	$(Y - \bar{Y})^2$
78	84	-5.7	-4.4	25.08	32.49	19.36
77	82	-6.7	-6.4	42.88	44.89	40.96
85	82	1.3	-6.4	-8.32	1.69	40.96
88	85	4.3	-3.4	-14.62	18.49	11.56
87	89	3.3		1.98	10.89	0.36
82	90	-1.7		-2.72	2.89	2.56
81	88	-2.7	-0.4	1.08	7.29	0.16
77	92	-6.7	3.6	-24.12	44.89	12.96
76	83	-7.7	-5.4	41.58	59.29	29.16
83	89	-0.7	0.4	-0.42	.49	0.36
91	98	18.3	9.4	177.68	176.89	92.16
93	99	9.3	10.4	98.58	86.49	112.36
1004	1061			<u>288.66</u>	<u>486.68</u>	<u>362.92</u>

Eqn to the regression line of Y on X is

$$(Y - \bar{Y}) = b_{YX}(X - \bar{X})$$

$$b_{YX} = Y \cdot \frac{\sigma_Y}{\sigma_X} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{\sum (X - \bar{X})^2} = \frac{288.66}{486.68} = 0.593$$

Regression line of Y on X is

$$(Y - 88.4) = (0.593)(X - 83.7)$$

$$Y = (0.593)X + (88.4 - 49.634)$$

$$Y = (0.593)X + 38.7$$

Eqn to the regression line of X on Y is

$$(X - \bar{X}) = b_{XY}(Y - \bar{Y})$$

$$b_{XY} = Y \cdot \frac{\sigma_X}{\sigma_Y} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{\sum (Y - \bar{Y})^2} = \frac{288.66}{362.92} = 0.795$$

Regression line of X on Y is

$$(X - 83.7) = (0.795)(Y - 88.4)$$

$$X = (0.795)Y + (83.7 - 70.278)$$

$$X = (0.795)Y + 13.422$$

→ Estimate the price index of cotton corresponding to price index of wool when

$$= 94, \quad = ?$$

$$x = (0.795)y + 13.422$$

$$x = (0.795)(94) + 13.422$$

$$x = 74.73 + 13.422$$

$$x = 88.152$$

Curve Fitting

Method of Least Squares

$$y = f(x)$$

Normal equation

$$a_0 x^0 + a_1 x^1 + a_2 x^2 + \dots + a_n x^n, a_n \neq 0$$

$$y = a + bx$$

$$y_i = f(x_i) ; 1 \leq x_i \leq n$$

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

$$S = \sum_{i=1}^n e_i^2 = \sum \{y_i - f(x_i)\}^2$$

Fitting of a Straight line

$$y = a + bx \longrightarrow I$$

Normal eqn are

$$\sum_i y_i = na + b \sum x_i$$

$$\sum_i y_i x_i = a \sum x_i + b \sum x_i^2$$

→ Fit a suitable straight line to the following data.

x	1	2	3	4	5
y	14	21	40	55	68

x	y	xy	x^2
1	14	14	1
2	21	54	4
3	40	120	9
4	55	220	16
5	68	340	25
15	204	748	55

$$n=5, \sum x_i = 15, \sum y_i = 204, \sum x_i y_i = 748, \sum x_i^2 = 55$$

Normal eqn are

$$204 = 5a + 15b \rightarrow ①$$

$$748 = 15a + 55b \rightarrow ②$$

$$3 \times ① \Rightarrow 15a + 45b = 612$$

$$\begin{aligned} ② \Rightarrow & \cancel{15a + 55b = 748} \\ & +10b = +136 \end{aligned}$$

$$b = 13.6$$

$$\text{From eq } ①, 5a + 15(13.6) = 204$$

$$5a = 0$$

$$a = 0$$

Required eqn of the straight line is

$$= (13.6)x$$

Fitting of a power curve

$$y = ab^x$$

$$\log y = \log(ab^x)$$

$$\log y = \log a + \log b^x$$

$$\log y = \log a + x \log b$$

$$Y = A + Bx$$

$$\text{where } A = \log a$$

$$B = \log b$$

$$Y = \log y$$

Normal equations are

$$\sum_i Y_i = nA +$$

$$\sum_i x_i Y_i = A \sum_i x_i + B \sum_i x_i^2$$

→ Fit a power curve to the following data.

x	0	1	2	3	4	5	6	7
y	10	21	35	59	92	200	400	600

x	y	$Y = \log_{10} y$	xy	x^2
0	10	1.0000	0.0000	0
1	21	1.3222	1.3222	1
2	35	1.5440	3.0880	4
3	59	1.7708	5.3124	9
4	92	1.9638	7.8552	16
5	200	2.3010	11.5050	25
6	400	2.6020	15.6120	36
7	600	2.7853	19.4971	49
28		15.2891	64.1919	140

$$n=8, \sum_i x_i = 28, \sum_i y_i = 15.2891, \sum_i x_i y_i = 64.1919, \sum_i x_i^2 = 140$$

Normal eqn are

$$\sum_i y_i = nA + B \sum_i x_i$$

$$\sum_i x_i y_i = A \sum_i x_i + B \sum_i x_i^2$$

$$15.2891 = 8A + 28B \rightarrow ①$$

$$64.1919 = 28A + 140B$$

$$2.2926 = A + 5B \rightarrow ②$$

Solving ① & ②

$$① \Rightarrow 15.2891 = 8A + 28B$$

$$\begin{array}{r} ② \times 8 \Rightarrow 12.3 = 8A + 40B \\ - \\ \hline +12B = +3.051 \end{array}$$

$$B = \frac{3.051}{12}$$

$$B = 0.2543$$

From eq ②,

$$A + 5$$

3. Random Variables

Axioms

$$0 \leq P(E) \leq 1$$

$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

$$P(\Omega) = 1$$

$$P(E) + P(\bar{E}) = 1$$

Addition Rule

$$P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2)$$

where $P(E_1 \cap E_2)$ = Joint probability

Conditional Probability

$$P(E_1 | E_2) = \frac{P(E_1 \cap E_2)}{P(E_2)} ; P(E_2) \neq 0$$

$$P(E_2 | E_1) = \frac{P(E_1 \cap E_2)}{P(E_1)} ; P(E_1) \neq 0$$

Independent events

$$P(E_1 | E_2) = P(E_1)$$

$$P(E_2 | E_1) = P(E_2)$$

Multiplication Rule/ Compound theorem on Probability

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2 | E_1)$$

$$= P(E_2) \cdot P(E_1 | E_2)$$

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2) \quad (\text{for independent events})$$

→ Find the probability of getting a sum 10 or 11 when two dice are thrown

$$E_1 \rightarrow 10$$

$$E_2 \rightarrow 11$$

die-I	die-II
6	4
5	5
4	6

die-I	die-II
6	5
5	6

$$P(E_1) = \frac{3}{36}$$

$$P(E_2) = \frac{2}{36}$$

$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

$$= \frac{3}{36} + \frac{2}{36}$$

$$= \frac{5}{36}$$

→ A bag contains 4 white & 6 black balls. Two balls are drawn ^{out} ~~at~~ random in succession. What is the probability that the two balls drawn out are white,

- (i) the first ball drawn is replaced in the bag before the second ball is drawn.
- (ii) the first ball drawn is not replaced in the bag before the second ball is drawn.

$$E_1 \rightarrow \text{first ball drawn}$$

$$E_2 \rightarrow \text{second ball drawn}$$

$$P(E_1 \cap E_2) = ?$$

$$P(E_1) = \frac{4}{10}$$

(i) with replacement

$$P(E_2) = \frac{4}{10}$$

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2)$$

$$= \frac{4}{10} \times \frac{4}{10}$$

(ii) without replacement

$$P(E_2 | E_1) = \frac{3}{9}$$

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2 | E_1)$$

$$= \frac{4}{10} \times \frac{3}{9}$$

Discrete Probability Distribution

A discrete random variable X assumes x_1, x_2, \dots, x_n with the corresponding probabilities $P_1, P_2, P_3, \dots, P_n$, then the following tabular form is known as

Discrete Probability Distribution.

$$X : x_1 \ x_2 \ x_3 \ \dots \ x_n$$

$$P(X=x_i) : P_1 \ P_2 \ P_3 \ \dots \ P_n$$

$$P_i = P(X=x_i) \quad 1 \leq i \leq n$$

where P_i = probability mass function (pmf)

↓

It must satisfy the following conditions

$$(i) \ P_i \geq 0,$$

$$(ii) \ \sum_i P_i = 1$$

→ Consider a discrete random variable X assumes the no. of heads turned up when two coins are tossed simultaneously.

$$X : S \rightarrow R$$

$$x(HH) = 2$$

$$S = \{S_1, S_2, S_3, S_4\}$$

$$x(HT) = 1$$

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

$$x(TH) = 1$$

$$x(TT) = 0$$

$$X : 0 \ 1 \ 2$$

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

$$P(X=0) = P(TT)$$

$$= P(T \cap T)$$

$$= P(T) \cdot P(T)$$

$$= \frac{1}{2} \cdot \frac{1}{2}$$

$$= \frac{1}{4}$$

$$P(X=1) = P(HT \text{ or } TH)$$

$$= P(HT) + P(TH)$$

$$= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2}$$

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

$$P(X=2) = P(HH)$$

$$= P(H \cap H)$$

$$= P(H) \cdot P(H)$$

$$= \frac{1}{2} \cdot \frac{1}{2}$$

$$= \frac{1}{4}$$

Expectation of a Random Variable

$X \rightarrow$ discrete random variable

$$E(X) = \sum_i x_i P(X=x_i)$$

$$E(X^r) = \sum_i x_i^r P(X=x_i) \quad \text{where } r=1, 2, 3, \dots$$

$$E(X^2) = \sum_i x_i^2 P(X=x_i)$$

Mean and variance of discrete random variable

$$\text{Mean}(\mu) = E(X) = \sum_i x_i P(X=x_i)$$

$$\text{Variance}(\sigma^2) = E(X-\mu)^2 = E(X^2)-\mu^2$$

Properties:

(i) $E(K) = K$, where 'K' is any constant

(ii) $E(X+Y) = E(X) + E(Y)$

(iii) $E(aX+bY) = aE(X) + bE(Y)$

(iv) $E(aX+b) = aE(X) + b$

(v) $V(K) = 0$; where V = Variance

(vi) $V(aX+b) = a^2 V(X)$

(vii) $E(XY) = E(X) \cdot E(Y)$; if X, Y are independent

→ For the following discrete probability distribution

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

Find (i) K

(ii) $P(X < 6)$

(iii) $P(X \geq 6)$

(iv) $P(0 < X < 5)$

(v) Mean and variance

X is a discrete random variable

$$\sum_{i=1}^6 p_i = 1$$

$$0 + K + 2K + 2K + 3K + K^2 + 2K^2 + 7K + K = 1$$

$$10K^2 + 9K = 1$$

$$10K^2 + 9K - 1 = 0$$

$$10K^2 + 10K - K - 1 = 0$$

$$10K(K+1) - (K+1) = 0$$

$$(K+1)(10K-1) = 0$$

$$K = -1 \text{ (not)} \quad K = \frac{1}{10}$$

$$\because p_i \geq 0, \forall i, \text{ then } K = \frac{1}{10} \checkmark$$

$$P(X \leq 6) = P(X = 0 \text{ or } 1 \text{ or } 2 \text{ or } 3 \text{ or } 4 \text{ or } 5)$$

$$= P(0) + P(1) + P(2) + P(3) + P(4) + P(5)$$

$$= 0 + K + 2K + 2K + 3K + K^2$$

$$= K^2 + 8K$$

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

$$= \frac{81}{100}$$

$$= 0.81$$

$$P(X \geq 6) = P(X = 6 \text{ or } 7)$$

$$= P(6) + P(7)$$

$$= 2K^2 + 7K^2 + K$$

$$= 9K^2 + K$$

$$= 9\left(\frac{1}{100}\right) + \frac{1}{10}$$

$$= \frac{19}{100}$$

$$= 0.19$$

$$P(0 < X < 5) = P(X = 1 \text{ or } 2 \text{ or } 3 \text{ or } 4)$$

$$= K + 2K + 2K + 3K$$

$$= 8K$$

$$= 8\left(\frac{1}{10}\right)$$

$$= 0.8$$

$$\text{Mean}(\mu) = E(x)$$

$$= \sum x_i P(x=x_i)$$

$$\mu = 0(0) + 1(k) + 2(2k) + 3(2k) + 4(3k) + 5(k^2) + 6(2k^2) + 7(1k^2+k)$$

$$\mu = 6k^2 + 30k$$

$$\mu = \frac{66}{100} + \frac{30}{10}$$

$$\mu = \frac{366}{100}$$

$$\boxed{\mu = 3.66}$$

$$\text{Variance}(\sigma^2) = E(x^2) - \mu^2$$

$$E(x^2) = \sum x_i^2 P(x=x_i) = 0^2(0) + 1^2(k) + 2^2(2k) + 3^2(2k) + 4^2(3k) + 5^2(k^2) + 6^2(2k^2) + 7^2(1k^2+k)$$

$$= 440k^2 + 124k$$

$$= \frac{440}{100} + \frac{124}{10}$$

$$= 16.8$$

$$V(x) = E(x^2) - \mu^2$$

$$= 16.8 - (3.66)^2$$

$$= 3.4044$$

→ For the following discrete probability distribution.

x	-3	-2	-1	0	1	2	3
$P(x=x_i)$	k	0.1	k	0.2	$2k$	0.4	$2k$

Find (i) k

(ii) Mean

(iii) Variance

Continuous Probability Distribution

$$P\left(x - \frac{\partial x}{2} \leq X \leq x + \frac{\partial x}{2}\right) = f(x) dx$$

$$\frac{d}{dx}(F(x)) = f(x)$$

$f(x) \rightarrow$ probability density function (pdf)

$$f(x) \geq 0, \forall x$$

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

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

$$E(X^2) = \int_{-\infty}^{\infty} x^2 f(x) dx$$

$$E(X^r) = \int_{-\infty}^{\infty} x^r f(x) dx \text{ where } r=1, 2, 3, \dots$$

$$\mu = E(X)$$

$$V(X) = E(X^2) - \mu^2$$

$$\xrightarrow{①} f(x) = \begin{cases} K(1-x^2), & 0 < x < 1 \\ 0, & \text{otherwise} \end{cases}$$