

Unit - I

* Discrete and continuous rate :- Discrete data is the information than can take only certain values.

Eg:- shoes size

These type of data is often represented using tally charts, bar charts or pie charts.

Continuous data is the data that can take any value

Eg:- Speed, height, weight, temperature.

This data is shown on, line graph, histogram & scatter plot

* Measure of central tendency :- One of the most important objective of statistical analysis is to get one single value that describes the characteristics of the entire mass of data such a value is called a central value or average

* Objectives :- 1) Measures of central value by condensing the mass of data into one single value gives the wide view of entire data.

2) Measure of central value by reducing the marks data to one single figure enables comparisons to be made

* Requisites of good average :-

- 1) Easy to understand
- 2) Should be simple to compute so that it can be used widely
- 3) It should be based on all the items
- 4) should be well defined
- 5) should be capable for further treatment
- 6) It should have sampling stability.

* Types of central tendency :-

1) Arithmetic mean :- Arithmetic mean of the set of observations is their sum by the no of observations. If $x_1, x_2, x_3, \dots, x_n$ given n observations then the mean is given by $\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$

problem) The monthly income 12 families is given below
 280 180 96 98 104 85 80 94 100 75 600 200
 find the arithmetic mean

$$\Rightarrow \bar{x} = \frac{\sum x_i}{n} = \frac{1992}{12} = 166 \text{/-}$$

In case of freq. distribution

$$\begin{array}{l} x_1, x_2, x_3, \dots, x_n \\ f_1, f_2, f_3, \dots, f_n \end{array} \quad \bar{x} = \frac{f_1 x_1 + f_2 x_2 + \dots + f_n x_n}{f_1 + f_2 + \dots + f_n}$$

$$= \left[\bar{x} = \frac{\sum f_i x_i}{\sum f_i} \right] \quad [x = \sum f]$$

2) find the arithmetic mean of

$$\begin{array}{ccccccc} x: & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ f: & 5 & 9 & 12 & 17 & 14 & 10 & 6 \end{array}$$

$$\Rightarrow \bar{x} = \frac{5+18+36+68+70+60+42}{73} = 4.095 \text{/-}$$

3) calculate mean marks for the foll. data

$$\begin{array}{cccccc} \text{Marks:} & 0-10 & 10-20 & 20-30 & 30-40 & 40-50 & 50-60 \\ \text{No. of stu:} & 12 & 18 & 27 & 20 & 17 & 6 \end{array}$$

Marks	No. of stu (f)	Mid pt (x)
0-10	12	5
10-20	18	15
20-30	27	25
30-40	20	35
40-50	17	45
50-60	6	55

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{60+270+675+700+765+530}{100} = 28 \text{/-}$$

$$\bar{x} = \frac{2800}{100} = 28 \text{/-}$$

* Step deviation method :- when n & f values are large the previous used formulae become time consuming hence in case of group or continuous freq. distribution the operation can be reduced to a greater extent by taking $d = \frac{x-A}{h}$ where A - assumed mean / arbitrary value
 h - common magnitude of class interval & arithmetic mean is given by

$$\bar{x} = A + \frac{\sum f_i d}{\sum f_i} \times h$$

C.I.	freq	Mid. pt (x)	f.x
0-8	8	4	32
8-16	7	12	84
16-24	16	20	320
24-32	24	28	672
32-40	18	36	576
40-48	7	44	280

$$\Rightarrow \bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{32+84+320+672+576+280}{80} = \frac{1964}{80} = 24.55 \text{/-}$$

C.I.	f	Mid. pt (x)	$d = \frac{x-A}{h}$	f.d	Assumed mean $(A) = 28$
0-8	8	4	-3	-24	
8-16	7	12	-2	-14	
16-24	16	20	-1	-16	
24-32	24	28	0	0	$h = 8$
32-40	18	36	1	18	
40-48	7	44	2	14	
				80	
					-22

$$\bar{x} = A + \frac{\sum f_i d \times h}{N} = 28 + \frac{(-22) \times 8}{80} = 25.8 \text{/-}$$

$$\bar{x} = 28 - \frac{176}{80}$$

$$\bar{x} = 25.8 \text{/-}$$

* Properties of mean :- Algebraic sum of the deviation of the set of values from their arithmetic mean is zero.

2) If \bar{x}_i for $i=1, 2, \dots, k$ are the means of k component series of sizes n_i for $i=1, 2, \dots, k$ then the mean \bar{x} is obtained by combining the component series to obtain the formula

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + n_3 \bar{x}_3 + \dots + n_k \bar{x}_k}{n_1 + n_2 + n_3 + \dots}$$

q:-> The mean weight of 100 workers in a factory running two shifts of 60 & 40 workers with the mean of 38.2, the mean age of 60 workers is 240 find the mean age of 40 workers in the night shift.

$$\Rightarrow n_1 = 60, n_2 = 40$$

$$\bar{x} = 38$$

$$\bar{x}_1 = 40, \bar{x}_2 = 38$$

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2} \quad \bar{x} = \frac{60(40) + 40(38)}{60 + 40} \quad \bar{x}(100) = 2400 + \bar{x}(40)$$

$$38 \times 100 = 2400 + 40 \bar{x}_2 \quad 1400 = 40(\bar{x}_2) \quad \therefore \bar{x}_2 = 1400 / 40$$

2) The avg weekly salary of male employee in a firm was 5200 & that of female of 4200 the mean salary of all the emp was 5000. find the percentage of male & female employee

$$\Rightarrow \bar{x}_1 = 5200, \bar{x}_2 = 4200, \bar{x} = 5000$$

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}; 5000 = \frac{n_1(5200) + n_2(4200)}{n_1 + n_2}$$

$$5000 n_1 + 5000 n_2 = 5200 n_1 + 4200 n_2 \quad \frac{n_1}{n_2} = \frac{800}{200} = 4$$

$$800 n_2 = 200 n_1$$

$$n_1 = \frac{4}{4+1} \times 100 = 80 \therefore, n_2 = \frac{1}{4+1} \times 100 = 20 \therefore$$

* Median :- Median is the middle item of the series arranged in ascending or descending order of magnitude
(*) if the no. of observation is odd the median is the middle value.

(ii) If the no. of observations are even then median is obtained by taking arithmetic mean of middle terms.

$$\text{Ex:- } 17 \ 25 \ 20 \ 15 \ 45 \ 18 \ 78 \ 12 \ 38 \ 19 \\ \Rightarrow 10 \ 12 \ 15 \ 18 \ 20 \ 25 \ 38 \ 45 \ 78 \ 19 \Rightarrow \text{median} = 28 \therefore$$

$$2) 8 \ 20 \ 50 \ 25 \ 15 \ 30$$

$$\Rightarrow 8 \ 15 \ 20 \ 25 \ 30 \ 50 \quad \text{median} = 22.5 \therefore$$

b) For continuous frequency distribution :- The median is obtained by using formula:-

$$\text{Median} = l + \frac{h}{2} \left[\frac{N}{2} - c \right] \quad \text{where } l \rightarrow \text{lower limit of median class}$$

f - frequency of the median class, h - magnitude of median class
c - Its cumulative freq. of the classes preceding the median class

$$N = 5 \therefore$$

1) find the median of the foll. distribution of class interval

C.I	2000 - 3000	3000 - 4000	4000 - 5000	5000 - 6000	6000 - 7000
f	3	5	20	10	5

$$\Rightarrow M/2 = 5f/2 = 43/2 = 21.5$$

C.I	f	C.f
2000 - 3000	3	3
3000 - 4000	5	8
(4000 - 5000)	20	28
5000 - 6000	10	38
6000 - 7000	5	43

$$\text{median} = 4000 - 5000$$

$$l = 4000, h = 1000, f = 20, C = 8$$

$$\text{median} = l + \frac{h}{f} \left[\frac{N}{2} - c \right]$$

$$= 4000 + \frac{1000}{20} [21.5 - 8] = 4675 \therefore$$

2) find the median of the foll.

C.I	0-7	7-14	14-21	21-28	28-35	35-42	42-49
f	19	25	36	22	51	43	28

$$\Rightarrow \frac{\sum f}{2} = \frac{224}{2} = 112$$

C.I	f	C.F
0-7	19	19
7-14	25	44
14-21	36	80
21-28	22	102
28-35	51	153
35-42	43	196
42-49	28	224

$$\text{median class} = 28 - 35$$

$$l = 28, h = 7, f = 51, C = 8$$

$$\begin{aligned}\text{Median} &= l + \frac{h}{f} \left[\frac{N}{2} - C \right] \\ &= 28 + \frac{7}{51} \left[\frac{112}{2} - 8 \right]\end{aligned}$$

$$\text{median} = 42.27\text{,}$$

3) In a factory employing 3000 persons in a day 5% work less than 3 hrs, 580 work from 3.01 to 4.50 hrs, 30% work from 4.51 to 6.00 hrs, 500 work from 6.01 to 7.50 hrs 20% work from 7.51 to 9.00 hrs & the rest work 9.01 or more what is the median hrs at work.

C.I	f	C.F	class boundaries
less than 3	150	150	less than 3.005
3.01-4.50	580	730	3.005 - 4.505
4.51-6.00	900	1630	4.505 - 6.005
6.01-7.50	500	2130	6.005 - 7.505
7.51-9.00	600	2730	7.505 - 9.005
more 9.01	270	3000	more than 9.005

$$\frac{\sum f}{2} = \frac{5f}{2} = \frac{3000}{2} = 1500$$

$$\text{median} = 4.505 + 1.5 \left(\frac{1500 - 27}{900} \right)$$

$$\text{median class} = 4.505 - 6.005$$

$$\text{median} = l + \frac{h}{f} \left[\frac{\sum f}{2} - C \right]$$

$$l = 4.505, h = 1.5, f = 900$$

$$C = 2730$$

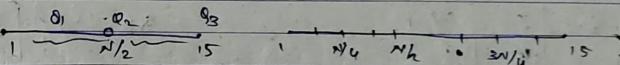
$$\text{median} = 5.788\text{,}$$

* Measure of dispersion :- consider scores of 2 cricketers A & B. A = {34, 6, 23, 81, 76, 9, 127} ; B = {28, 41, 36, 33, 38, 27, 35} for batsman A the score differ 34 (mean) and differ large but batsman B the score differ from 34 q diff is small just based on this we can decide who is more consistent i.e. the property which denotes the extent to which the individual values of the variables are dispersed about central tendency is called dispersion

* Quartile Deviations (Q.D.) - It is also called as semi-interquartile deviation

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

$$\text{Median} = l + \frac{h}{f} \left[\frac{N}{2} - C \right]$$



$$Q_1 = l + \frac{h}{f} \left[\frac{N}{4} - C \right]$$

$$Q_3 = l + \frac{h}{f} \left[\frac{3N}{4} - C \right]$$

$$\text{Quartile Deviation (Q.D.)} = \frac{Q_3 - Q_1}{3}$$

$$\text{where } Q_1 \rightarrow l + \frac{h}{f} \left[\frac{N}{4} - C \right] \quad Q_3 \rightarrow l + \frac{h}{f} \left[\frac{3N}{4} - C \right]$$

This gives only the absolute measure of this version but for the comparative study of variability b/w two distribution relative measure known as coefficient of quartile range

$$\text{coeff. of Q.D.} = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

1) calculate the quartile deviation & coeff. quartile deviation for the foll. ungrouped items. 20 28 40 12 30 15 50

$\Rightarrow 10 \ 15 \ 20 \ 25 \ 30 \ 40 \ 50$

$$Q_1 = \text{size of } \left(\frac{n+1}{4}\right)^{\text{th}} \text{ item}$$

$= \text{size of } 2^{\text{nd}} \text{ item}$

$= 15$

$$Q_3 = \text{size of } \left(\frac{3(n+1)}{4}\right)^{\text{th}} \text{ item}$$

$= \text{size of } 6^{\text{th}} \text{ item}$

$= 40$

$$Q.D = \frac{Q_3 - Q_1}{2} = \frac{40 - 15}{2} = \frac{25}{2} = 12.5,$$

$$\text{Coeff of Q.D} = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{40 - 15}{40 + 15} = \frac{25}{55} = \frac{5}{11} = 0.45$$

2) calculate semi-interquartile range & coeff of the full data

C.I	0-10	10-20	20-30	30-40	40-50	50-60	60-70	70-80
f	11	18	25	28	30	33	22	15

C.I	f	C.f
0-10	11	11
10-20	18	29
20-30	25	54
30-40	28	82
40-50	30	112
50-60	33	145
60-70	22	167
70-80	15	182

$$Q_1 = l + \frac{h}{f} \left[\frac{n}{4} - c \right]$$

$$Q_3 = l + \frac{h}{f} \left[\frac{3n}{4} - c \right]$$

for Q_1

$$\frac{n}{4} = \frac{112}{4} = 28$$

median class = 20-30, $f = 25$,

$l = 10$, $h = 10$, $c = 29$.

$$\therefore Q_1 = 20 + \frac{10}{25} [45.5 - 29]$$

$$Q_1 = 26.6$$

for $Q_3 \Rightarrow \frac{3n}{4} \Rightarrow 136.5$

median class = 50-60, $f = 33$, $h = 10$, $l = 50$, $c = 112$

$$Q_3 = 20 + \frac{10}{33} [136.5 - 112] \quad Q_3 = 57.42$$

$$Q.D = Q_3 - Q_1 \Rightarrow \frac{57.42 - 26.6}{2} \Rightarrow 15.41$$

$$\text{Coeff of Q.D} = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{57.42 - 26.6}{57.42 + 26.6} = 0.3668$$

Now

3) find Q.D for data

C.I	40-49	50-59	60-69	70-79	80-89	90-99
f	306	192	144	96	42	34

$$\Rightarrow Q.D \approx 10.31, \text{ Coeff of Q.D} = 0.189.$$

* Standard Deviation :- this gives the +ve sq. root of the arithmetic mean of the sq's of deviations of given observations from their Arithmetic mean.

$$S.D = \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}} \quad \bar{x} = \sqrt{\frac{\sum x_i^2 - (\sum x)^2}{N}}$$

7) calculate the standard deviation the ungrouped data

x	x - \bar{x}	(x - $\bar{x})^2$
8	-4	16
9	-3	9
15	3	9
23	11	121
5	-7	49
11	-1	1
19	7	49
8	-4	16
10	-2	4
12	0	0

$\bar{x} = \frac{\sum x}{N}$

$$S.D = \sqrt{\frac{\sum (x - \bar{x})^2}{N}} = \sqrt{\frac{274}{10}} = 5.2311$$

2) Given the full information find SD $n=10$, $\sum x = 60$, $\sum x^2 = 1000$

$$\sigma = \sqrt{\frac{\sum x^2}{N} - \left(\frac{\sum x}{N}\right)^2} = \sqrt{\frac{1000}{10} - \left(\frac{60}{10}\right)^2} = \sqrt{100 - 36} = 8.11$$

2) Calculate the S.D. for the foll. data.

size of items	6	7	8	9	10	11	12
freq	3	6	9	12	8	5	4

C.I	f	$x - \bar{x}$	$(x - \bar{x})^2$	$f(x - \bar{x})^2$
6	3	-3	9	27
7	6	-2	4	24
8	9	-1	1	9
9	13	0	0	0
10	8	1	1	8
11	5	2	4	20
12	4	3	9	36
			124	

$$\bar{x} = \frac{18 + 42 + 72 + 117 + 80 + 55 + 48}{48} = 9$$

$$\therefore \sigma = \sqrt{\frac{\sum (x - \bar{x})^2 f}{\sum f}} = \sqrt{\frac{124}{48}} = 10.607$$

3) find the S.D. for foll. data.

C.I	0-10	10-20	20-30	30-40	40-50	50-60	60-70	70-80
f	15	15	23	22	25	10	5	10

C.I	f	mid pt (m)	$x - \bar{x}$	$(x - \bar{x})^2 f$
0-10	15	5	-30.16	13644.384
10-20	15	15	-20.16	6096.384
20-30	23	25	-10.16	2374.184
30-40	22	35	0.16	0.5632
40-50	25	45	9.16	2420.64
50-60	10	55	19.16	3936.256
60-70	5	65	29.16	4452.125
70-80	10	75	39.16	15822.25

$$49801.78$$

$$\bar{x} = \frac{\sum xf}{\sum f} = \frac{35 + 225 + 575 + 770 + 1125 + 550 + 325 + 70}{48} = 35.16$$

6) $\sigma = \sqrt{\frac{\sum f(x - \bar{x})^2}{\sum f}}$

$$\sigma = \sqrt{\frac{49801.78}{125}} = 19.4583$$

3) the analysis of the result of a budget survey of 150 families gave an average monthly expenditure of £ 120 on food items with SD of £ 15. After the analysis completed it was noted that fig. recorded for the household was wrongly taken as £ 15 instead of £ 105 determine correct value of avg expenditure & its std. deviation

$$\Rightarrow n = 150, \bar{x} = 120, \text{ wrong value } \sigma = 15, \text{ correct value } = 105$$

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

$$120 = \frac{\sum x}{150}, \sum x = 18000$$

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

$$15 = \sqrt{\frac{\sum x_i^2 - (120)^2}{150}}, 225 + 14400 = \frac{\sum x_i^2}{150}$$

$$\sum x_i^2 = 2193750$$

$$\sum x = 18000 - 15 \times 105 = 18090$$

$$\bar{x} = \frac{\sum x}{n} = \frac{18090}{150} = 120.6$$

$$\sum x_i^2 = 2193750 - 15^2 + 105^2 = 2204550$$

$$\sigma = \sqrt{\frac{\sum x_i^2 - (\bar{x})^2}{n}}, \sigma = \sqrt{\frac{2204550}{150} - \frac{(120.6)^2}{150}}$$

$$\sigma = \sqrt{8977} = 9.47$$

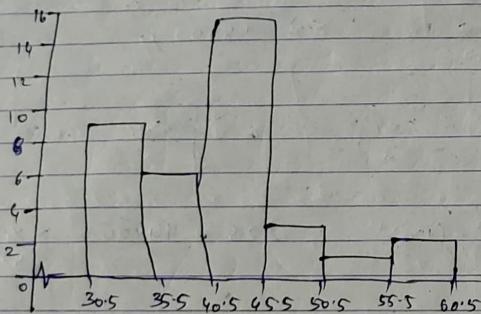
* Graphical Representation of Data :-

* Histogram :-

1) Draw the Histogram for the following data

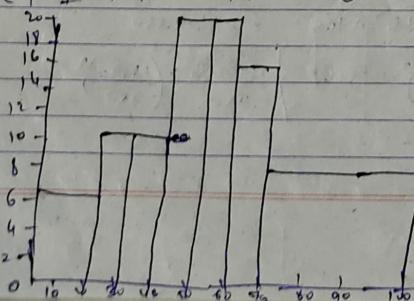
weight (kg)	No. of students
30.5 - 35.5	9
35.5 - 40.5	16
40.5 - 45.5	15
45.5 - 50.5	3
50.5 - 55.5	1
55.5 - 60.5	2

⇒



2) A teacher wanted to analyze the performance of 2 sections of students in a maths test of 100 marks. The following marks were formed.

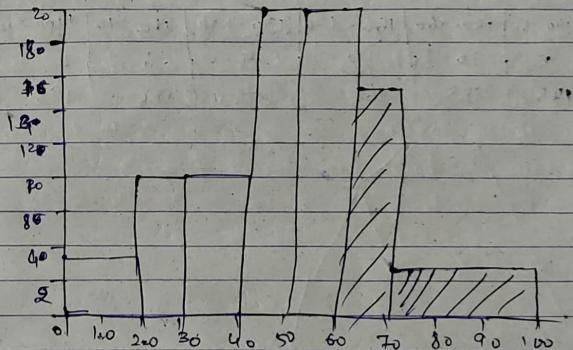
Marks	0 - 20	20 - 30	30 - 40	40 - 50	50 - 60	60 - 70	70 & above
No. of St.	7	10	10	20	20	15	8



Since the intervals are not same here

(2) select a class interval of min class size here the class size is 10, the length of the rectangle are modified to be proportionate to the class size.

Marks	No. of students (f)	width of class	length of rectangle
0 - 20	7	20	$\frac{7}{10} \times 10 = 7$
20 - 30	10	10	$\frac{10}{10} \times 10 = 10$
30 - 40	10	10	$\frac{10}{10} \times 10 = 10$
40 - 50	20	10	$\frac{20}{10} \times 10 = 20$
50 - 60	20	10	$\frac{20}{10} \times 10 = 20$
60 - 70	15	10	$\frac{15}{10} \times 10 = 15$
70 & above	8	30	$\frac{8}{10} \times 10 = 8$

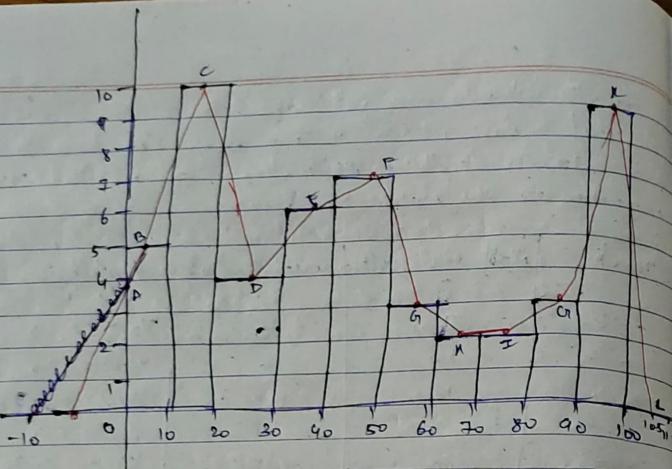


* Frequency polygon :-

q^2 → Draw the freq. polygon for the foll. data.

Marks	0 - 10	10 - 20	20 - 30	30 - 40	40 - 50	50 - 60	60 - 70	70 - 80	80 - 90	90 - 100
No. of St.	5	10	4	6	7	3	2	2	3	9

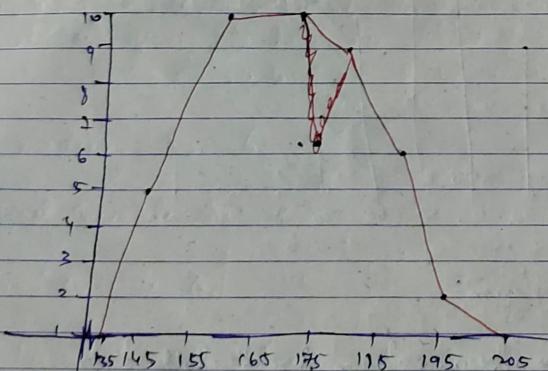
⇒



Ques draw freq-poly for 1st problem of histogram

2) draw freq-poly without using histogram

Cost of living.	160-150	150-160	160-170	170-180	180-190	190-200
no. weeks	5	10	10	9	6	2
mid pt	165	155	165	175	185	195



* Stem leaf graph :-

1) Draw the stem leaf graph for the foll data.

15 27 8 17 13 22 24 25 30 36 38 32 37

28 43 7

steam	leaf
1 3	7 8
2 2	3 3 5 7
3 2	2 4 5 7 8
4 3	2 2 6 8
6	3

2) Draw the stem leaf graph for the foll data.

1.2, 2.3, 2.5, 2.4, 2.6, 1.8, 2.7, 3.2, 4.1, 2.9, 4.5,
2.6, 5.8, 9.3, 10.6, 12.4, 10.9

⇒

steam	leaf
1	0.2 0.8
2	0.3 0.4 0.5 0.6
3	0.2
4	0.1 0.5
5	0.8
7	0.6
9	0.3
10	0.6 0.9
12	0.4

* Box plot

1) draw the box plot of the foll data 10, 12, 11, 15, 11, 14, 13, 17, 12, 12, 14, 11

⇒ 10, 11, 11, 14, 12, 12, 13, 14, 14, 15, 17, 22

$$\text{Median} = Q_2 = 12.5$$

$$\text{First quartile} = Q_1 = \frac{11+11}{2} = 11$$

$$\text{Third quartile} = Q_3 = \frac{14+15}{2} = 14.5$$

min value = 10 max value = 20 outlier is above

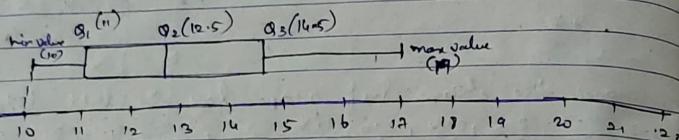
(3 types - less than, more than, Both combination).

$$\text{Interquartile range (IQR)} = Q_3 - Q_1 = 14.5 - 11 = 3.5$$

$$\text{Range} [Q_1 - (1.5 \times IQR), Q_3 + (1.5 \times IQR)]$$

$$[11 - (1.5 \times 3.5), 14.5 + (1.5 \times 3.5)]$$

$$[5.75, 19.75]$$



min value = 10, max value = 22, outlier is 22

The owner of the restaurant wants to find more about where his workers are coming from. He decided together with his workers to gather about distances that people commuted to the rest. People reported the foll. dist. travelled.

14, 6, 3, 2, 4, 15, 11, 8, 1, 7, 2, 1, 3, 4, 10, 15, 20
Create a graph to help him spread of dist's

$\Rightarrow 1, 1, 2, 2, 3, 3, 4, 4, 6, 7, 8, 10, 11, 12, 14, 15, 20$

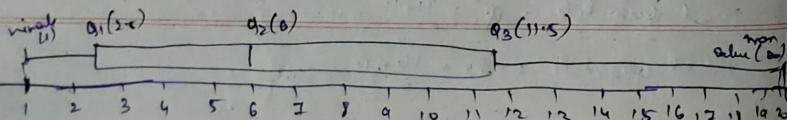
$$\text{median} = Q_2 = 6$$

$$\text{First Quartile } Q_1 = 2.5$$

$$\text{Third Quartile } Q_3 = 11.5$$

$$\text{Interquartile range (IQR)} = Q_3 - Q_1 = 11.5 - 2.5 = 9.0$$

$$\text{Range} [Q_1 - (1.5 \times 9), Q_3 + (1.5 \times 9)] = [-11, 25]$$



* Ogive curves :-

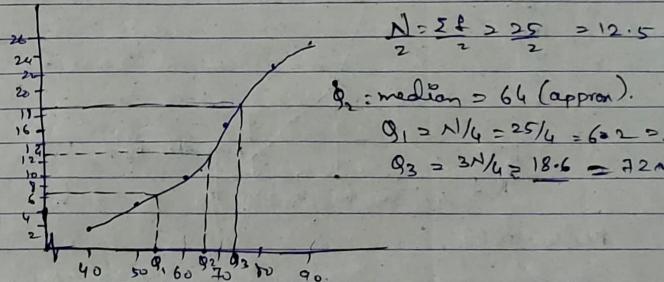
a) plot the Ogive curve for the foll. data.

C.I	30-40	40-50	50-60	60-70	70-80	80-90
f	2	3	5	7	6	2

\Rightarrow

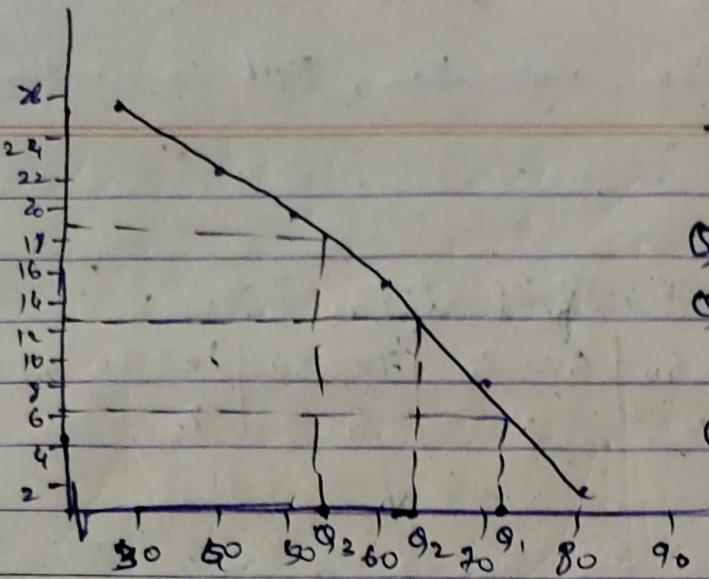
b) less than curve

C.I	f	less than	c.f
30-40	2	40	2
40-50	3	50	5
50-60	5	60	10
60-70	7	70	17
70-80	6	80	23
80-90	2	90	25



b) More than curve

C.I	f	more than	c.f
30-40	2	30	25
40-50	3	40	23
50-60	5	50	20
60-70	7	60	15
70-80	6	70	8
80-90	2	80	2



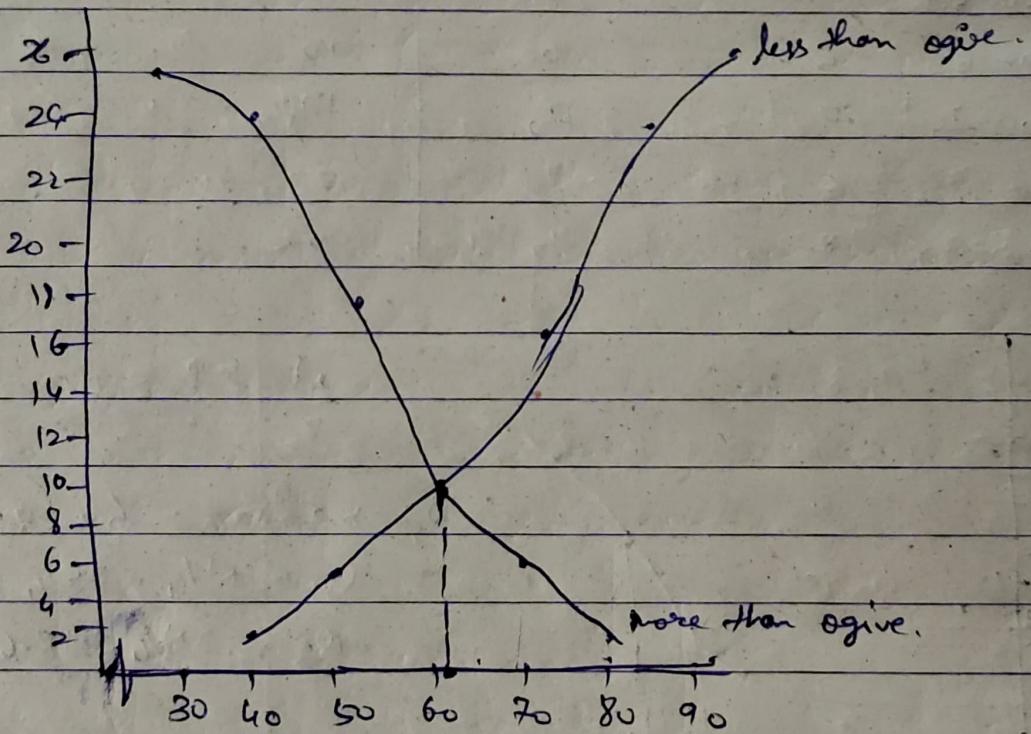
$$\frac{N}{2} > \frac{\sum f}{2} > \frac{25}{2} = 12.5$$

$Q_2 = \text{median} = 62 \text{ ~n}$

$$Q_1 = \frac{N}{4} = \frac{25}{4} = 6.25 = 71 \text{ ~n}$$

$$Q_3 = \frac{3N}{4} = \frac{75}{4} = 18.75 = 52 \text{ ~n}$$

(iii) Both combination



Median ≈ 64