

A quantity which can vary from one individual to another is called variable or variate

PRIMARY DATA: original & first hand information
SECONDARY DATA: copy & second hand collected from records or already available

GEOGRAPHICAL CLASSIFICATION
→ Data is classified as per geographical region or place.
→ Eg. Production of wheat in diff. region. countries.

QUALITATIVE CLASSIFICATION
→ Data is classified on the basis of quality.

QUANTITATIVE CLASSIFICATION AND FREQUENCY DISTRIBUTION

Let the marks obtained by 40 students of AI out of 50 marks

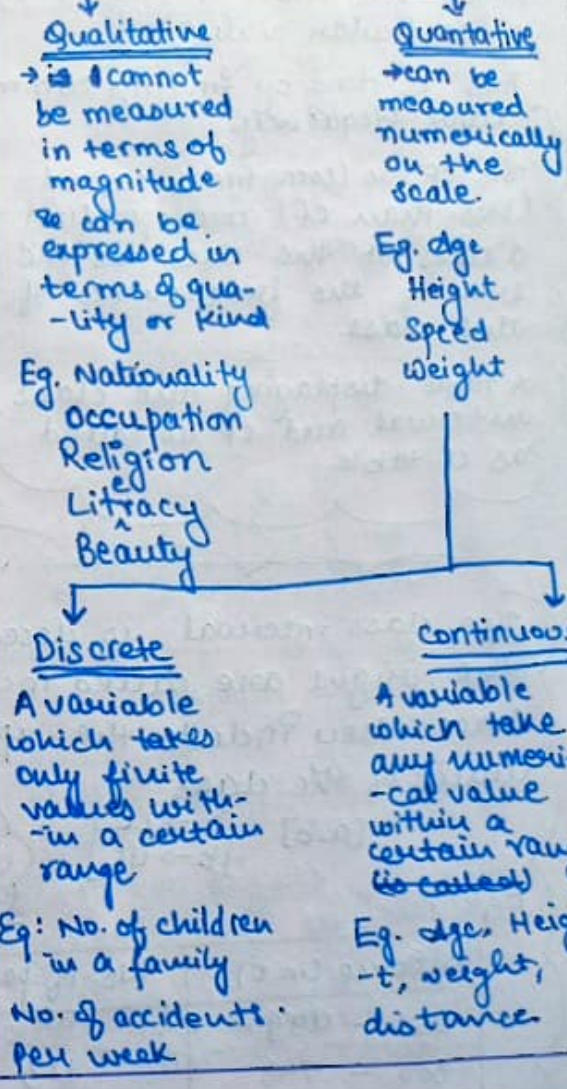
5, 6, 7, 8, 11, 15, 20, 8, 11, 25, 30, 15, 17, 11,
6, 22, 25, 20, 22, 15, 30, 32, 32, 8, 20, 25, 22,
22, 35, 37, 40, 20, 11, 25, 10, 20, 10, 25, 35, 42

Marks	Sample Space	Tally Bar	Frequency	CF
5-10	5		7	40
10-15	5		6	33
15-20	5		5	27
20-25	5		9	22
25-30	5		4	13
30-35	5		4	9
35-40	5		3	5
40-45	5		2	2

greater than
Sample space must be equal in all groups

CF
7
7+6
7+6+5
7+6+5+9
7+6+5+9+4
...

TYPES OF VARIABLE



Summary of dec-4

- Definition of C.F
- Concept of Exclusive & Inclusive class interval
- How to convert Inclusive class into Exclusive class interval

CUMULATIVE FREQUENCY

Some times one may want to know less than or greater than a particular value.

This is done by finding cumulative frequency.

The CF is (sometimes called less than CF) corresponding to a class is the sum of the frequencies of that class.

A table displaying the class interval and CF is called as CF table

The class interval is closed to the right are called inclusive since they include the upper limit of the class

$[a, b]$ $\begin{cases} a \rightarrow LL \\ b \rightarrow UL \end{cases}$

EQ:

Income (in ₹)	No. of person
250 - 499	60
500 - 749	43
750 - 999	25

Compare class limits before adj & after adj

However, for further statistical analysis is desirable that class interval is exclusive.

To convert inclusive class int. into ECI: ~~one is~~

- ① Find the diff. b/w LL of second class to the UL of the first class.
- ② Subtract the half of this diff. from all the ~~the~~ LL and add to all the UL.

- ③ The adjustment factor in the above eg. is: $\frac{500 - 499}{2} = 0.5$

- ④ Now, from all the lower limits \rightarrow add 0.5 and from all the upper limit \rightarrow subtract 0.5

EXCLUSIVE AND INCLUSIVE CLASS INTERVALS

The class intervals open to the right are called Exclusive as they exclude the upper limit of the class

$[a, b)$ $\begin{cases} a \rightarrow \text{lower limit} \\ b \rightarrow \text{upper limit} \end{cases}$

EQ:

Income (in ₹)	No. of person
250 - 500	75
500 - 750	70
750 - 1000	52
1000 - 1250	33

\rightarrow exclusive class interval
In ECI, the upper limit of upper class is the lower limit of next class

If the income of a person is 750 then he is included in class 750 - 1000.

Lecture 5 Summary

16/01/2024

- Histogram formation
- Frequency Polygon formation

GRAPHICAL REPRESENTATION OF DATA

The grouped frequency distribution provides a better idea of the data as compared to the ungrouped frequency dist.

Further, if the distribution is represented by graphs a more clear visual impression about the data is obtained as graphs are the ~~grouped~~ good visual aid.

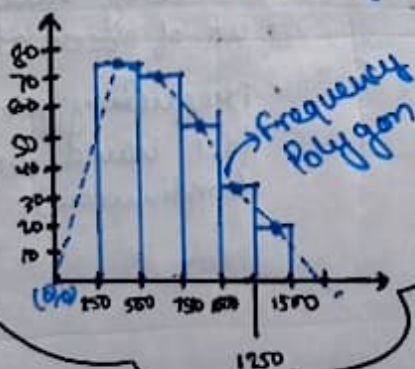
various imp. types of graphs :-

- Histogram
- Bar Charts
- Frequency curve
- Frequency polygon

② FREQUENCY POLYGON (FP)

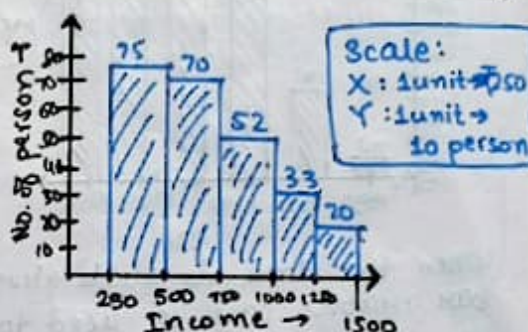
For a grouped frequency distribution with equal class interval a FP is obtained by:-

- ① Joining the middle points of the top of the rectangles of the histogram by means of dotted straight lines.
- ② To complete the polygon the mid pt. at each end are joined to the just lower & higher mid pts of zero frequency (origin)



① HISTOGRAM

INCOME (in ₹)	No. of Persons
250 - 500	75
500 - 750	70
750 - 1000	52
1000 - 1250	33
1250 - 1500	20



To draw the histogram

- ① Mark all the class interval on X-axis
- ② Mark frequencies along Y-axis
- ③ It is not necessary to take same scale on both axis. You can take different scales.
- ④ Construct rectangles for each class interval having the height equal to the corresponding frequency.

• Bar Chart

- Simple Bar Diagram
- Multiple Bar Diagram

→ Frequency Curve

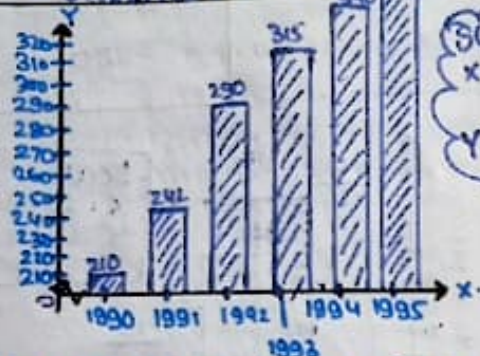
BAR CHART

These are also called as bar graphs. Various types of bar graphs:-

→ Simple Bar Diagram

E.g.

YEAR	1990	1991	1992	1993	1994	1995
No. of STUDENTS	210	242	290	315	340	355



When the data is qualitative or categorical bar graphs can be used to represent the data.

(i) A bar graph can be drawn horizontally or vertically.

NOTE: The height of the bar represent the frequency of data.

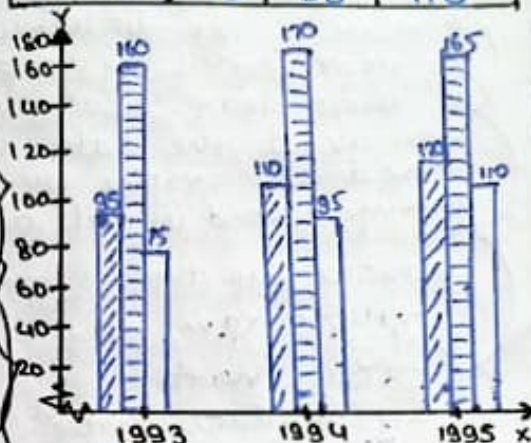
(ii) The distance b/w two bars should be constant

(iii) The width of each bar must be same

→ Multiple Bar Diagram

E.g.

YEARS	1993	1994	1995
ARTS	95	110	120
SCIENCE	160	170	165
COMMERCE	75	95	110



SCALE

x-axis: 1 unit → 1 year

y-axis: 1 unit → 20 students

= Arts

= Science

= Commerce

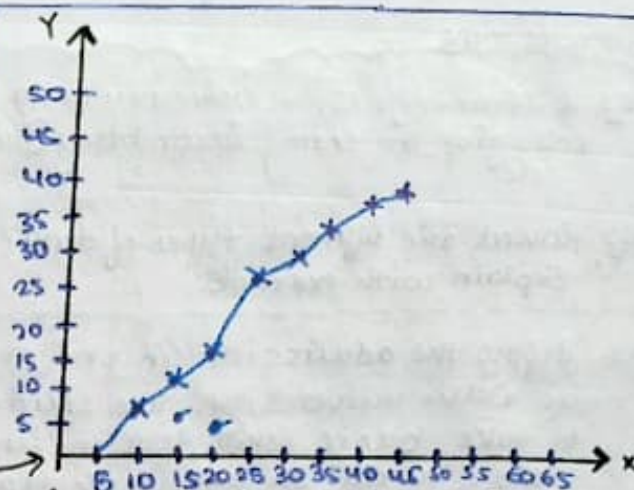
→ Frequency Curve

- Grouped frequency distribution
- width of class int. → very small
- no. of observation → increased
- Frequency polygon → free hand smooth curve (continuous)

Lecture-7 (Summary)

- less than ogive curve plotting
- more than ogive curve plotting

C.I	Tally	#freq	CF (less than)	CF (greater than)
5-10		7	7	40
10-15		6	13	33
15-20		5	18	27
20-25		9	27	22
25-30		4	31	13
30-35		4	35	9
35-40		3	38	5
40-45		2	40	2



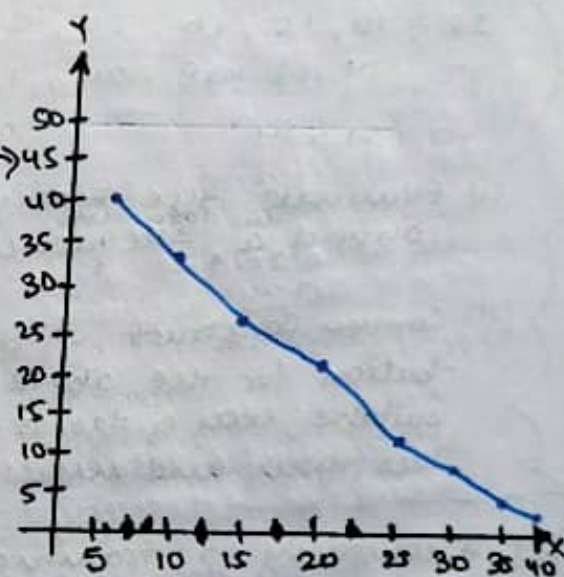
less than ogive

There are two types of curve

- less than ogive
- more than ogive

For less than ogive:-

→ we plot the points corresponding to upper limit on X-axis and (less than) CF on Y-axis. A free hand smooth curve and joining the points we get less than ogive.



More than ogive

For more than ogive:-

→ we plot the points corresponding to lower limits on X-axis and (more than) CF on Y-axis. A free hand smooth curve joining the points we get more than ogive

IMP QUES TYPE

SCALE-MEASUREMENT OF SCALE

Q1. Difference b/w - Primary Data & Secondary Data (also give an exam. -ple)

Methods for collecting Primary Data & Sec. Data

Q2. What are different types of data? Explain with example

Q3. Before the admission in phd program in Delhi university. The student have to take basic skill test in fundamen-tals of math. In one such example, 30 student appeared in the exam. They are marked out of 30 points.

15, 12, 15, 22, 28, 30, 19, 25, 24,
28, 10, 15, 16, 20, 26, 22, 18, 20,
27, 14, 12, 19, 21, 18, 19, 30, 13,
10, 21, 24

(i) Eliminate the steps you take to convert ^{this data into} a frequency distribution

(ii) ~~Convert~~ Construct a frequency distri-bution for the above data with a suitable no. of classes. Also construct less than and more than ogive.

Q4. A company a training program after the first week, the company officers joining the training program. The score out of 100 of 30 employees:-

32, 36, 31, 67, 65, 42, 39, 56, 78, 61,
34, 78, 75, 78, 61, 30, 65, 45, 48, 78
43, 75, 64, 73, 87, 41, 56, 71, 81, 85

(i) convert the data into frequency distribution by taking a suitable width.

(ii) Construct a histogram & frequency poly-gon diagram.

- Scale
- Mean
- Measures of central Tendency

SCALE

1. Nominal Scale: on basis of quality
2. Ordinal Scale: on basis of order
3. Interval Scale: " " " interval

MEASURE OF CENTRAL TENDENCY

- Single value that attempt to describe a set of data by identifying the central position within that set of data.
- As such of central tendency are sometimes called measures of central location
- It is a way of summarising the data in a form of a typical or ^{representative} logical values
- Four types of mean:-
 - Arithmetic mean
 - Geometric mean
 - Harmonic mean

CENTRAL TENDENCY OF DATA

1. Mean

$$\rightarrow \text{Mean} = \frac{\sum x_i f_i}{\sum f_i}$$

$$\rightarrow \sum (x - m) = 0$$

$$\rightarrow \sum (x - m)^2 \rightarrow \text{minimum}$$

$$\rightarrow \text{mean} = \mu$$

$$\text{Mean of sample} = \bar{X}$$

$$\rightarrow \bar{X} = \frac{\sum fx}{N} \quad [\text{Mean of group ed Data}]$$

→ use mean when:-

- data are interval or ratio scale
- data are not skewed
- because it is sensitive to every score

23rd Jan '24

Lecture -10 Summary

- Basic formulae
- Mean \rightarrow merit \rightarrow Direct Method
- Median \rightarrow ways to calculate \rightarrow Step deviation method
- Median \rightarrow merit & Demerit

$$\text{Variance} = \frac{\sum (x_i - \bar{x}_i)^2}{n} \quad [\text{in form of mean}]$$

$$\text{Coeff. of variance} = \frac{\text{Mean}}{\text{S.D.}}$$

$$\text{Variance} = \frac{\sum (x_i - \text{Median})^2}{n} \quad [\text{in form of median}]$$

Q) $x: 3.2 \quad 5.8 \quad 7.9 \quad 4.5$
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $f_i: x \quad x+2 \quad x-3 \quad x+6$
 Find x if $AM = 4.876$

$$A = \bar{x} = \frac{\sum x_i f_i}{\sum f_i}$$

Substitute & get the answer

$$4.876(4x+15) = 2.143 + 14.9x$$

$$x = 30.717$$

MEDIAN

\rightarrow Positional number which divides the given series into two parts.

ANSWER: Median is 3rd position

MERITS

- Not affected by extreme obs.
- can be computed with open & closed intervals.

DEMERIT:

- In case of even number of obs. for grouped data, median cannot be determined exactly
- Not based on each and every items of distribution
- cannot be used to compute the median of combined group
- Median is less stable because it depends on its position.

Mean

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i}$$

Arithmetic Mean

DEMERIT:

- 1. AM is affected by extreme observations
- 2. AM can not be used in case of open & closed classes.
- 3. AM cannot be used to calculate when dealing with qualitative information
- 4. AM can be used if any one observation is missing.

DIRECT METHOD

$$x = \frac{L+U}{2} \quad [\text{mid pt. of CI}]$$

$$\bar{x} = \frac{\sum x_i f_i}{\sum f_i}$$

STEP DEVIATION

$$\text{Mean} = A + h \frac{\sum f d}{\sum f}$$

$$d = \frac{x - A}{h}$$

\downarrow Middle value of CI $\quad \downarrow$ assumed mean $\quad \downarrow$ class size

Marks	No. of St.	Mid point	$d = \frac{x-A}{h}$	$f d$
0-10	6	5	-3	-18
10-20	5	15	-2	-10
20-30	8	25	-1	-8
30-40	29	35	0	0
40-50	7	45	1	7
50-60	6	55	2	12
60-70	3	65	3	9
	$\sum f = 60$			$\sum f d = -8$

$$h = 10$$

$$A = 35$$

$$x = \frac{L+U}{2}$$

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

$$\bar{x} = 35 + 10 \frac{-8}{60}$$

$$\bar{x} = 33.4$$

January 2024

• Question Solving

lecture - 11 Summary

Question

Calculate the mean & std. deviation for the following table given age distribution of 542 members

Age Group	f
20 - 30	3
30 - 40	61
40 - 50	132
50 - 60	153
60 - 70	140
70 - 80	51
80 - 90	2

dim $h = 10$; $A = 55$

Age	f	X	$d = \frac{X-A}{h}$	fd	fx_i	fd^2
20-30	3	25	-3	-9	75	27
30-40	61	35	-2	-122	2135	246
40-50	132	45	-1	-132	5940	132
50-60	153	55	0	0	8415	0
60-70	140	65	1	140	9100	140
70-80	51	75	2	102	3825	204
80-90	2	85	3	6	170	18
Sum	542			-15	29660	765

variance formula

$$\sigma^2 = h^2 \left[\frac{1}{N} (\sum fd^2) - \left(\frac{1}{N} (\sum fd) \right)^2 \right]$$

$N = \sum f$

(iii) Variance =

$$100 \left[\frac{1}{542} (765) + \left(\frac{1}{542} \times -15 \right)^2 \right]$$

$$= 100 [1.411 + 0.0007]$$

$$= 141.17$$

$$\text{Standard Deviation} = \sqrt{\text{Variance}}$$

$$= \sqrt{141.17}$$

$$= 11.88 \text{ Years}$$

(i) Direct Method -

$$\bar{X} = \frac{\sum fix_i}{\sum fi} = \frac{29660}{542} = 54.72$$

(ii) Step deviation

$$\bar{X} = A + h \left(\frac{\sum fd}{\sum f} \right)$$

$$= 55 + 10 \left(\frac{-15}{542} \right)$$

$$= 55 - 0.27 = 54.73$$

31st January 2024

Lecture - 12 Summary

- Question on mean
- Median
- Question on median of continuous class series

① Calculate the average marks of the following students

MARKS	NO. OF STUDENTS
0-10	5
10-20	12
20-30	15
30-40	25
40-50	8
50-60	3
60-70	2

C.I	f	X	$d = \frac{X-A}{h}$	fd	fd^2	fix_i
0-10	5	5	-3	-15	45	25
10-20	12	15	-2	-24	48	180
20-30	15	25	-1	-15	15	375
30-40	25	35	0	0	0	875
40-50	8	45	1	8	8	360
50-60	3	55	2	6	12	165
60-70	2	65	3	6	18	130
	$\Sigma f = 70$			$\Sigma fd = -34$	$\Sigma fd^2 = 146$	$\Sigma fix_i = 2110$

MEDIAN

(i) For discrete series

$$\left(\frac{n+1}{2}\right)^{th}; \text{ when 'n' is odd}$$

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

(ii) For continuous series

$$\text{Median} = l + \frac{\frac{N}{2} - cf}{f} \cdot h$$

l = lower limit of median class

h = class size of median class

f = frequency of median class

N = Total freq. Σf_i

cf = cumulative frequency just preceding the median class

Step Deviation

1) Direct Method:-

$$\bar{X} = A + h \left(\frac{\Sigma fd}{\Sigma f} \right)$$

$$= 35 + 10 \left(\frac{-34}{70} \right)$$

$$= 30.14 \text{ Marks}$$

2) Direct Method

$$\bar{X} = \frac{\Sigma fix_i}{\Sigma f_i}$$

$$= \frac{2110}{70}$$

$$= 30.14 \text{ Marks}$$

3) Variance

$$\sigma^2 = h^2 \left[\frac{\Sigma fd^2}{N} - \left(\frac{\Sigma fd}{N} \right)^2 \right]$$

$$\sigma^2 = 100 \left[\frac{1}{70} \times 146 - \left(\frac{-34}{70} \right)^2 \right]$$

$$\sigma^2 = 100 [2.08 - 0.23]$$

$$\sigma^2 = 185$$

4) Step deviation (σ) = $\sqrt{\sigma^2}$

$$= \sqrt{185} = 13.6$$

④ Find the median of data

Age	10-15	15-20	20-25	25-30	30-35	35-40	40-45	45-50	50-55	55-60	60-65
Person	8	12	10	8	3	2	7				

$$f = 10; l = 20; h = 10; N = 50; cf = 20$$

$$\text{Median} = l + \frac{\frac{N}{2} - cf}{f} \cdot h$$

$$= 20 + \frac{25 - 20}{10} \times 10 = 20 + 5 = 25 \text{ day}$$

AGE	f	cf
0-10	8	8
10-20	12	20
20-30	10	30
30-40	8	38
40-50	3	41
50-60	2	43
60-70	7	50
	$\Sigma f = 50$	

$$\frac{N}{2} = \frac{50}{2} = 25$$

25 is just greater than cf.

To find median class

- Mode
- Question Practice

MODE

$$\text{Mode} = l + h \frac{(f_1 - f_0)}{(2f_1 - f_0 - f_2)} \quad (\text{for continuous case})$$

l = lower limit of modal class

h = class size of modal class

f_0 = frequency of class preceding the modal class

f_1 = freq. of modal class

f_2 = freq. of class just ^{one} exceeding the modal class

For verification of Mode

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

never -ve

~~is not~~

Find the mode of the following:-

AGE	PERSON
0-10	8
10-20	12
20-30	10
30-40	8
40-50	3
50-60	2

AGE	PERSON
0-10	8
10-20	12
20-30	10
30-40	8
40-50	3
50-60	2

Modal class = highest freq.

↓
10-20

$l = 10$

$h = 10$

$$\begin{aligned} \text{Mode} &= 10 + 10 \frac{(12 - 8)}{(2 \times 12) - 8 - 10} \\ &= 10 + 10 \frac{4}{24 - 18} \end{aligned}$$

$$\begin{aligned} &= 10 + \frac{40}{6} \\ &= 10 + 6.67 \\ &= 16.67 \end{aligned}$$

Mode can be more than 1 in number

Method of grouping

Question

From the following data, regarding weight of 60 student of a class. Find mode.

weight = 50 : 51 : 52 : 53 : 54 : 55 : 56 : 57 : 58 : 59 : 60
 No. of st. = 2 : 4 : 5 : 6 : 8 : 5 : 4 : 7 : 11 : 5 : 3

④ In column (4), ^{sum up} combine the frequencies three-by-three starting from top.

⑤ In column (5), combine three frequencies three-by-three starting from second.

① In column (1) write the original frequency

② In column (2), combine the frequencies two-by-two starting from top.

③ In column (3), combine the frequencies two-by-two starting from second class.

⑥ In column (6), combine the freq. three-by-three starting from third.

weight	(1) f	(2)	(3)	(4)	(5)	(6)
50	2	6	///	11	///	///
51	4		9		15	19
52	5	11	14	19	17	
53	6		9			22
54	8	13	18	8	19	
55	5					11
56	4	11	16	8	19	
57	7					11
58	11	16	8	19		
59	5				11	23
60	3	///	///	///		

only applied on irregular frequency data

irregular data

Kabhi increase kr tha hai aur Kabhi decrease kr tha hai

Column	Max. Freq.	class corresponding to max. freq.			
(1)	11			58	
(2)	16			58	59
(3)	18		57	58	
(4)	22	56	57	58	
(5)	23		57	58	59
(6)	19	52 53 54	57	58	59 → 3 60 → 1
			4	6	

If the frequencies are ascending then apply the normal formula of mode

$$\text{mode} = 58$$

The median & mode of the following salary distribution are 3350 & 3400 resp.
Sum of all freq. = 230.

SALARY	No. of EMPLOYEE
0-1000	4
1000-2000	16
2000-3000	F_3
3000-4000	F_4
4000-5000	F_5
5000-6000	6
6000-7000	4

SALARY	f	CF (less)
0-1000	4	4
1000-2000	16	20
2000-3000	F_3	$20 + F_3$
3000-4000	F_4	$20 + F_3 + F_4$
4000-5000	F_5	$20 + F_3 + F_4 + F_5$
5000-6000	6	$26 + F_3 + F_4 + F_5$
6000-7000	4	$30 + F_3 + F_4 + F_5$

Median class

Modal class

$$② N = 230 = 20 + F_3 + F_4 + F_5$$

$$\Rightarrow 2000 = F_3 + F_4 + F_5 \quad \text{--- (i)}$$

③ Median 3350 lies in b/w 3000-4000

④ mode

$$= l + h \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right)$$

$$= 3000 + 1000 \left(\frac{F_4 - F_3}{2F_4 - F_3 - F_5} \right)$$

$$3400 = 3000 + 1000 \left(\frac{F_4 - F_3}{2F_4 - F_3 - F_5} \right)$$

$$\frac{2}{8} = \frac{F_4 - F_3}{2F_4 - F_3 - F_5}$$

From (i) & (ii):
 $F_4 + 2.85F_3 = 271.42$
 $F_4 - 3F_3 = -2F_5$

$$5.85F_3 = 271.42 + 2F_5$$

$$271.42 + 2F_5 = 5.85F_3$$

$$12F_3 + 2F_5 = 6000$$

$$1.5F_3 + F_5 = 728.58$$

$$10.85F_3 + 4F_5 = 4728.58$$

$$\text{Median} = l + \frac{h}{f} \left(\frac{N}{2} - cf \right)$$

$$3350 = 3000 + \frac{1000}{F_4} \left(\frac{230}{2} - 20 - F_3 \right)$$

$$350 = \frac{1000}{F_4} (115 - 20 - F_3)$$

$$F_4 = \frac{1000}{350} (95 - F_3)$$

$$F_4 = 2.85 (95 - F_3)$$

$$F_4 = 271.42 - 2.85F_3$$

$$F_4 + 2.85F_3 = 271.42 \quad \text{--- (i)}$$

From (i) & (iii):

$$2000 - F_3 - F_5 - 3F_3 = -2F_5$$

$$-4F_3 + F_5 = 2000 \quad \text{--- (ii)}$$

From (ii) & (iv):

$$F_4 - 3F_3 + 2F_5 = 0$$

$$F_4 - 1.85F_3 - F_5 = -1728.58$$

$$-1.85F_3 + 3F_5 = 1728.58 \quad \text{--- (vi)}$$

From (i) & (iv):

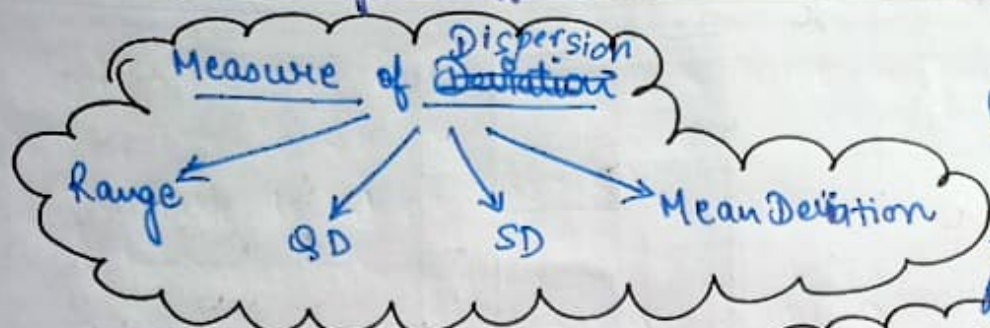
$$F_4 + 2.85F_3 = 271.42$$

$$-4F_3 + F_5 = 2000$$

$$F_4 - 1.85F_3 - F_5 = -1728.58$$

- Measure of Dispersion
- Range
- Range's coefficient

- Interquartile Range



$$COR = \frac{75-60}{75+60} \times 100$$

$$= \frac{15}{135} \times 100$$

$$= 11.11\%$$

Range

$$Range = U - L$$

upper limit (x_{max})

lower limit (x_{min})

calculate Range from the following distribution

SIZE:	60-63	63-66	66-69	69-72
NUM.:	5	1	42	27

$$\frac{72-75}{8}$$

$$Range = 75 - 60 = 15$$

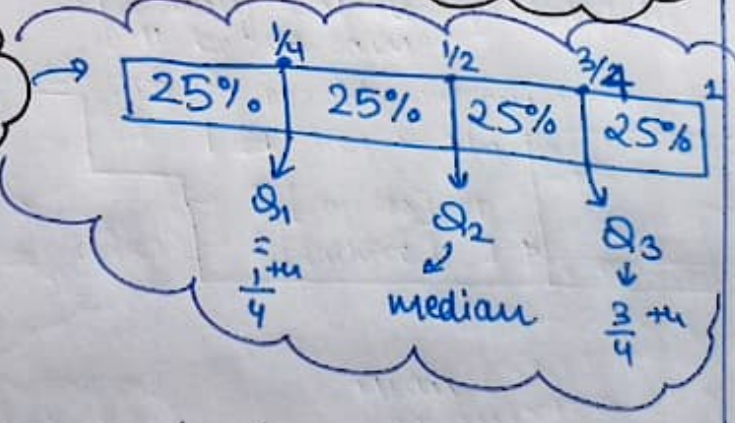
The yields of a cotton variety from 5 plots are 8, 9, 8, 19. Find the Range!

$$Range = 19 - 8 = 11$$

$$Coeff. of Range (COR) = \frac{U-L}{U+L} \times 100$$

INTERQUARTILE RANGE

→ Diff. b/w the upper & lower quartile values in the set of values

$$IQR = Q_3 - Q_1$$


9th February 2024

Lecture-17 Summary

- Question on interquartile range
- Quartile for continuous range

INTERQUARTILE RANGE

Q 2, 8, 5, 10, 12

A ① Arrange in ascending order

2, 5, 8, 10, 12

② for discrete data set

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

$$= \left(\frac{5+1}{4} \right)^{\text{th}} \text{ term}$$

$$= (1.5)^{\text{th}} \text{ term}$$

$$Q_1 = \frac{2+5}{2} = \frac{7}{2}$$

$$Q_3 = \left[\frac{3}{4}(n+1) \right]^{\text{th}} \text{ term}$$

$$= \left(\frac{3}{4} \times 6 \right)^{\text{th}} \text{ term}$$

$$= \frac{9}{2} = 4.5^{\text{th}} \text{ term}$$

$$= \frac{10+12}{2} = 11$$

count the
n = total
no. of
data given

Quartile for continuous case

$$\text{Median} = l + \frac{\frac{N}{2} - cf}{f}$$

$$Q_1 = l + \frac{\frac{N}{4} - cf}{f}$$

$$Q_3 = l + \frac{\frac{3}{4}N - cf}{f}$$

Q Find,

(i) Q_1 , (ii) Q_2 , (iii) Q_3

(iv) Range (v) Interquartile Range of the following nos.

12, 5, 22, 30, 7, 86, 15

14, 42, 53, 25, 65

A

①

5, 7, 12, 14, 15, 22,

25, 30, 36, 42, 53,

65

$$n = 12$$

② Q_1

$$= \left(\frac{n+1}{4} \right)^{\text{th}} \text{ term} = \frac{13}{4}^{\text{th}}$$

$$= (3.25)^{\text{th}}$$

यहाँ 0.25 हो या

0.50 हो या

0.75 हो। हम हमेशा

उन दो नों को Avg लेंगे।

$$= \frac{12+14}{2}$$

$$= 13$$

③ Q_2

$$= \left(\frac{n+1}{2} \right)^{\text{th}} = \frac{13}{2} = 6.5^{\text{th}}$$

$$= \frac{22+25}{2}$$

$$= 23.5$$

④ Q_3

$$= \frac{3}{4}(n+1) = \frac{13}{4} \times 3 = 9.75$$

$$= \frac{36+42}{2}$$

$$= \frac{68}{2}$$

$$= 34$$

⑤ IQR = $Q_3 - Q_1$

$$= 34 - 13$$

$$\text{IQR} = 21$$

⑥ Range = $65 - 5 = 60$

12 Feb 2024

Variance in discrete case
Quantile Deviation

Normal Distribution
Binomial Distribution

Lecture - 18 Summary
Variance & Mean
1 Marker Formulae

$$\text{Var}(x) = \frac{1}{n} \sum f_i (x_i - \bar{x})^2 \quad [\text{For discrete case}]$$

$$\text{Quantile deviation} = \frac{Q_3 - Q_1}{2}$$

$$\text{Coefficient of Quantile Deviation} = \frac{Q_3 - Q_1}{Q_3 + Q_1} \times 100$$

Question

The no. of vehicles sold by a major Maruti showroom in a day was recorded for 10 working days. The data is given as:

DAY	FREQ.
1	20
2	15
3	18
4	5
5	10
6	17
7	21
8	19
9	25
10	28

→ arrange as per freq. in ascending then find quantile range & C.O.Q.

NORMAL DISTRIBUTION

Total No. of possible = 2^x → no. of trial

If no. of trials is exceedingly large than the normal trials possibility then we will use Binomial Distribution

BINOMIAL DISTRIBUTION

when trials are independent & identical & distributive.

→ In every trial the probability will be same.

→ हमें दो possible outcomes होते हैं 'Yes' & 'No'

$$P[\text{Success}] = p$$

$$P[\text{Failure}] = 1 - p$$

Prob. Mass. & Func.
→ $P(X=x) = {}^n C_x p^x q^{n-x}$
 $X = P(\text{Success})$

$$\text{Var}\left(\frac{x}{2} + c\right) = ?$$

$$\& \text{Var}(x) = a$$

$$\text{Var}(x) = a$$

$$\text{Var}(x) = a$$

$$\text{Var}(x) = a$$

$$\text{Mean}(x+5) = ?$$

$$\text{Mean}(x) = 4$$

$$\text{If mean}(x) = 4$$

$$\text{then mean}(x+5) = 4+5 = 9$$

$$\text{Var}(x+c) = a?$$

$$\text{If Var}(x) = a \& x$$

$$\Rightarrow \text{Var}(x) + \text{Var}(c) = x$$

$$\text{Mean}(2x+c) \& \text{Me}$$

$$\text{Mean}(x) = a$$

$$\Rightarrow \text{Mean}(2x) + \text{Mean}(c)$$

$$\Rightarrow 2a + c$$

$$\text{Var}(2x+c) = ? \& \text{Var}(x) = a$$

$$\Rightarrow \text{Var}(2x) + \text{Var}(c)$$

$$\Rightarrow 4a + 0$$

Probability Mass function

$$\rightarrow P(X=x) = {}^n C_x p^x q^{n-x}$$

\downarrow Prob. of Success \downarrow Prob. of failure

Mean $(X) = np$

var $(X) = npq$

mean $(X) >$ var (X)
 \rightarrow only in binomial dist.

Ques

6 dice are thrown 725 times. How many times do you want at least three dice to show 5 or 6.

Ans

$$P[X \geq 3] \times 725$$

n = 6 dice

or

$$(1 - P[X < 3]) \times 725$$

$$1 - (P(X=0) + P(X=1) + P(X=2)) \times 725$$

NORMAL DISTRIBUTION

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \left(e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \right);$$

(Prob. Dist. Funct.)

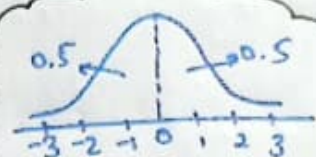
Mean (X) = $\mu = E(X)$

Var (X) = σ^2

$X \sim N(\mu, \sigma^2)$

condition

Also known as Gaussian Distribution



- Mean, Median, Mode are equal
- The curve is symmetric about mean.
- Total area under the curve is 1
- Left hand side & Right hand side of the graph about the mean are same.

Z-Score

$$Z = \frac{X - \mu}{\sigma}$$

Prob. Density ^{function.} ~~Distrib~~

$\mu = 0$

$\sigma = 1$

Question

If $Z = \frac{X - \mu}{\sigma}$ then $f(z) = ?$

Answer

$$f(z) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(z)^2}$$

• Questions on Normalisation

X is normally distributed & mean of X is 12 & std. deviation is 4. Find out the prob. of the following

a) $(X \geq 20)$

$$Z = \frac{X - \mu}{\sigma}$$

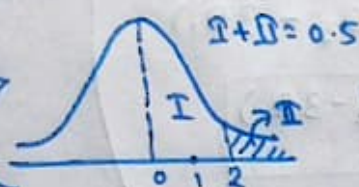
b) $X \leq 20$

c) $6 \leq X \leq 12$ $X \sim N(12, 16)$

$$P[X \geq 20] = P\left[\frac{X - \mu}{\sigma} \geq \frac{20 - \mu}{\sigma}\right]$$

↓

$$P[Z \geq 2] \leftarrow P\left[Z \geq \frac{20 - 12}{4}\right]$$



$$\rightarrow 0.5 - P[0 \leq Z \leq 2]$$

↓ given

$$\rightarrow 0.5 - 0.4772$$

$$\rightarrow 0.0228$$

Range of the graph

$$= -3 \text{ to } +3$$

Q Mean : 30

SD = 5

$X \sim N(30, 25)$

Find the prob.

(i) $26 \leq X \leq 40$

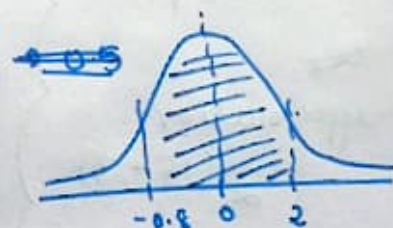
(ii) $X \geq 45$

(iii) $|X - 30| > 5$

$$P[26 \leq X \leq 40]$$

$$P\left[\frac{26 - \mu}{\sigma} \leq \frac{X - \mu}{\sigma} \leq \frac{40 - \mu}{\sigma}\right]$$

$$P\left[-\frac{4}{5} \leq Z \leq 2\right]$$



$$P[-0.8 \leq X \leq 0] + P[0 \leq Z \leq 2]$$

$$P[0 \leq X \leq 0.8] + P[0 \leq Z \leq 2]$$

$$0.2881 + 0.4772$$

$$= 0.7653 \text{ Ans}$$

coeff

$$\text{coefficient of variation} = \frac{\sigma}{\bar{x}} \times 100$$

σ = Std. Deviation

\bar{x} = Mean

To find the outliers:

$$\rightarrow \text{Lower outlier} = Q_1 - 1.5 IQR = Q_1 - 1.5 (Q_3 - Q_1)$$

$$\rightarrow \text{upper outlier} = Q_3 + 1.5 (IQR)$$

$$= Q_3 + 1.5 (Q_3 - Q_1)$$

$$= Q_3 (2.5) - (1.5 Q_1)$$

$$= \frac{1}{10} (25 Q_3 - 15 Q_1)$$

$$= \frac{15}{10} (5 Q_3 - 3 Q_1)$$

$$\text{upper outlier} = \frac{1}{2} (5 Q_3 - 3 Q_1)$$

$$= Q_1 - 1.5 Q_3 + Q_1 (1.5)$$

$$= +Q_1 (0.5) - 1.5 Q_3$$

$$= \frac{5}{10} (Q_1 - 3 Q_3)$$

$$\text{Lower outlier} = \frac{1}{2} (Q_1 - 3 Q_3)$$

Ques: 1, 3, 5, 8

$$\text{Ans: } Q_3 = \frac{3}{4} \times 4 = 3$$

~~upper~~

$$= 5$$

$$Q_1 = \frac{1}{4} \times 4 = 1$$

$$= 1$$

Q_3

upper outlier

$$= \frac{1}{2} (15 - 3)$$

$$= \frac{12}{2} = 6$$

upper outlier

$$= 5 + 1.5 (4)$$

$$= 5 + 6$$

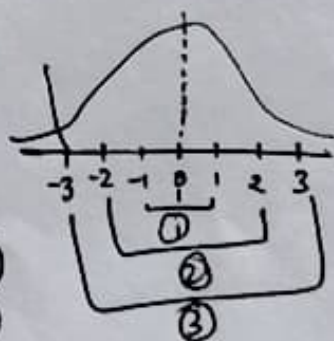
$$= 11$$

Area property of Normal dist

$$\textcircled{1} P [\mu - \sigma < x < \mu + \sigma] = 0.6826$$

$$\textcircled{2} P [\mu - 2\sigma < x < \mu + 2\sigma] = 0.9544$$

$$\textcircled{3} P [\mu - 3\sigma < x < \mu + 3\sigma] = 0.9973$$



20th February 2024

• Imp. Qⁿ/Topics

→ Types & Data & variables with real life eg. (long answer type)

→ Graphical representation of data (1-2 Qⁿ)

→ Central Tendency

→ Measures of Deviation (1-2 Qⁿ) → 100%

→ Outliers (1 Qⁿ short Ans. type)

→ Outliers using Interquartile range

→ Normal Distribution (using Z-score) for find probability.

COV

$$\frac{s}{\bar{x}} \times 100 = \text{COV}$$

$\text{var}(x+y) = \text{var}(x) + \text{var}(y) - 2\text{cov}(x,y)$
Note: If x & y are independent (diff.)
 then covar. of $(x \& y)$ are zero.

An analysis of monthly salary pay to the workers of 2 firms A & B belonging to the same industry gives the following results.

	<u>FIRMA</u>	<u>FIRM B</u>
No. of worker:	500	600
Avg. Salary:	186	175
variance of distribution	81	100
Std. Deviation:	9	10

→ Firm A
 $186 \times 500 = 93000$

→ Firm B
 $600 \times 175 = 105000$

(i) which firm A or B has a larger salary bill?

(ii) In which firm A or B is there greater variability in individual salary?

→ Firma
 $\text{COV} = \frac{9}{186} \times 100$
 $= 4.83$

→ Firm B
 $\text{COV} = \frac{10}{175} \times 100$
 $= 5.71$