

[Lecture Collected Power Point]

Statistics

Statistics is a branch of applied mathematics which is applied to observational data. It deals with data collection, presentation, organization and summarizing data, decision making etc.

Types:

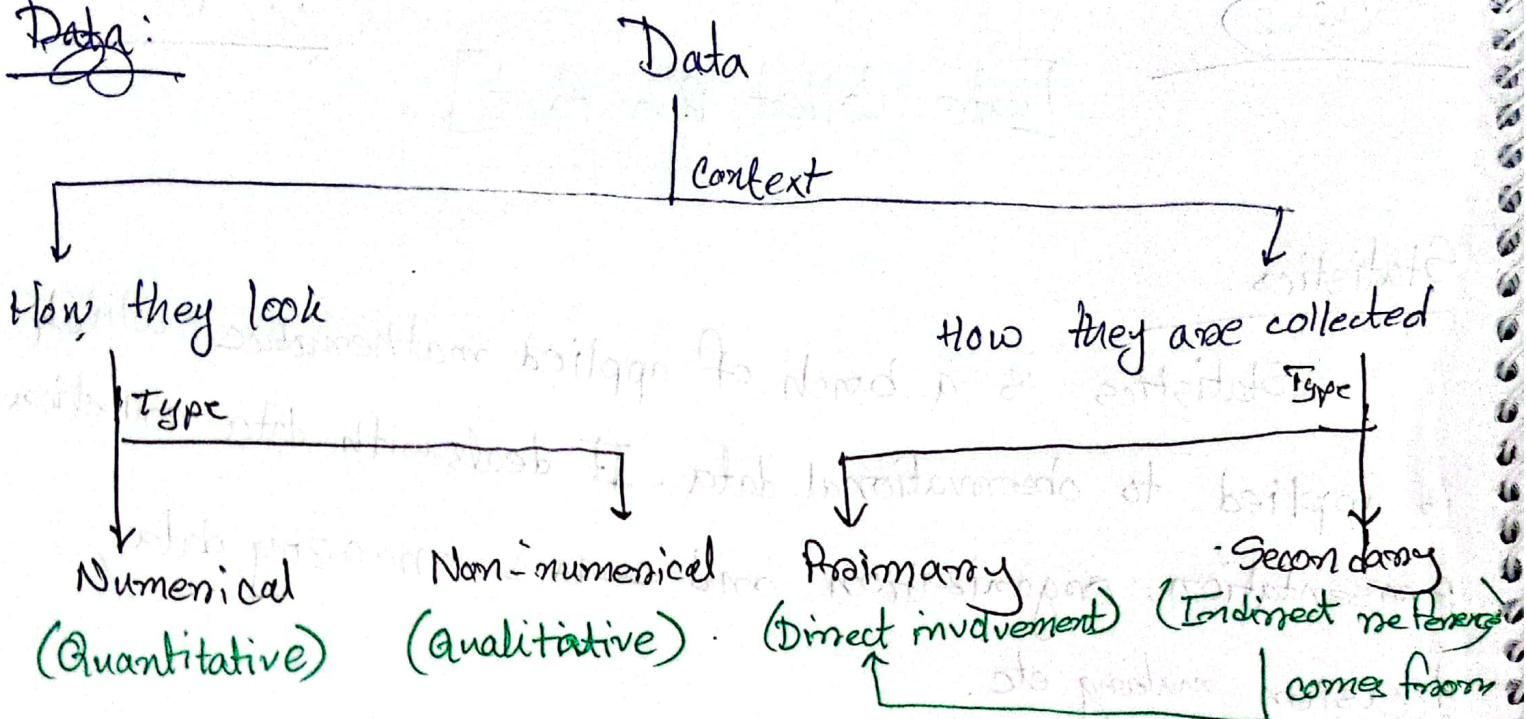
- i) Descriptive, [Everything upto Probability]
- ii) Inferential, [Af Theory of Probability TOP + F]
 - Theory of Probability, Sampling, Test of Significance, Test of hypothesis etc → based on sample → if random is involved
 - Data Presentation, Organization and summarization of data, measures of location, dispersion, shape characteristics of (central tendency things)

Population: The entire set of individuals or objects of interest on the measurements obtained from all individuals or objects of interest.

Sample: A $\xrightarrow{\text{representative}}$ portion, or part, of the population of interest

→ MUST Represent the population

Data:



Types of Levels of Data:

•

- i)
- ii) Interval Level Data → 0 means any non-existent
- iii) Ratio Level Data → 0 means non-existent and a state

• non-continuous but meaningful

• continuous & measurable

• ratio level data has a true zero point

• interval level data does not have a true zero point

• ratio level data has a true zero point

• interval level data does not have a true zero point

• ratio level data has a true zero point

• interval level data does not have a true zero point

III Describing Data:

* Frequency Distribution Table:

Construct a FDT / Arrange the data / Make a table.

$$k = \text{smallest number of classes}$$

$2^k \geq \text{no. of data}$: For what minimum value of k , 2^k will be greater than number of data

For 180,

$$2^8 = 256 > 180$$

Min class number = 256

i = class width, H = highest value, L = lowest value

$$i \geq \frac{H - L}{k}; \text{ always } [i]$$

* If gap is high, you can take b and adjust data

$$\left. \begin{array}{l} H = 3894 \\ L = 294 \\ i = 375 \end{array} \right\} \rightarrow \left. \begin{array}{l} H = 3900 \\ L = 200 \\ i = 400 \end{array} \right\}$$

Hence, the observation number is
 $n = 180$

[Check off from
Boole]

The smallest number of clauses (k) is,

$$2^k \geq 180$$

$2^8 \geq 180$ // find it in rough

which means, number of class, $k = 8$.

And, the class width (i) is,

$$i \geq \frac{H - l}{k} \Rightarrow \begin{cases} \text{Here,} \\ \text{highest value, } H = 3292 \\ \text{lowest } " ", l = 294 \end{cases}$$

$$\text{or, } i \geq \frac{3292 - 294}{8} \\ \geq 374.75$$

Hense,

We consider the class width 400
 We consider 200 in lieu of 294 and 300 in lieu of 3400 3400 of 3292
 e: Frequency Distribution Table of Profit (\$) on vehicles sold by

Table: Frequency Distribution Table of Profit (\$) on vehicles sold by

Applewood Auto // Name is MUST Write short but MUST

actual name
less than cumulative frequency

used to get a better idea of everything, proportion
Relative Frequency

class interval (P. of fit)	Tally	Frequency f_i	Cumulative frequency F_i	Relative Frequency $\left(\frac{f_i}{N}\right)$
200 - 600		8	8	$\frac{8}{180} = 0.04$
600 - - -	- -	11	19	$\frac{11}{180} = 0.06$
.	:	.	:	
.	:	1	:	
.	.	.	.	
$N = \sum f_i =$		180		

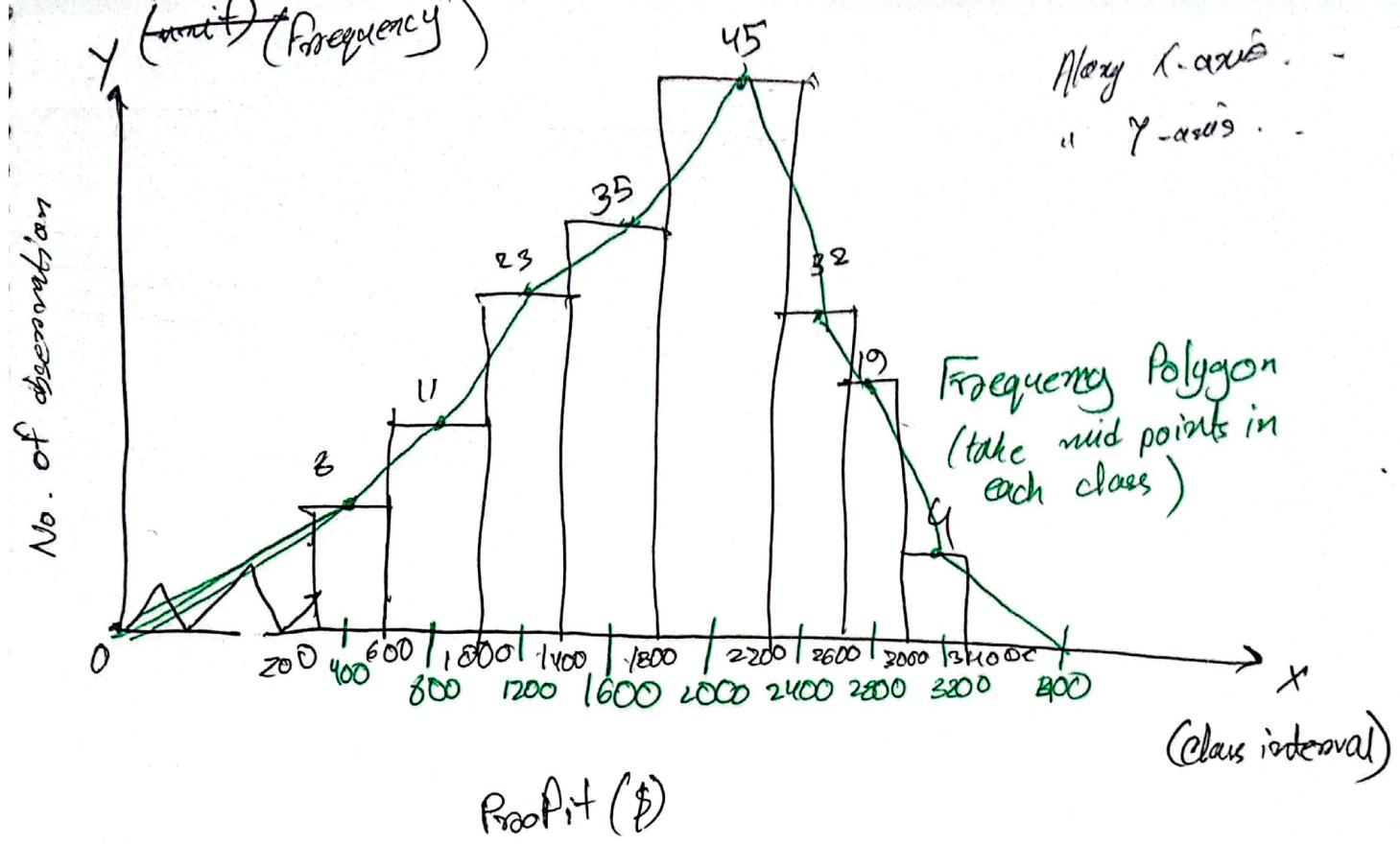
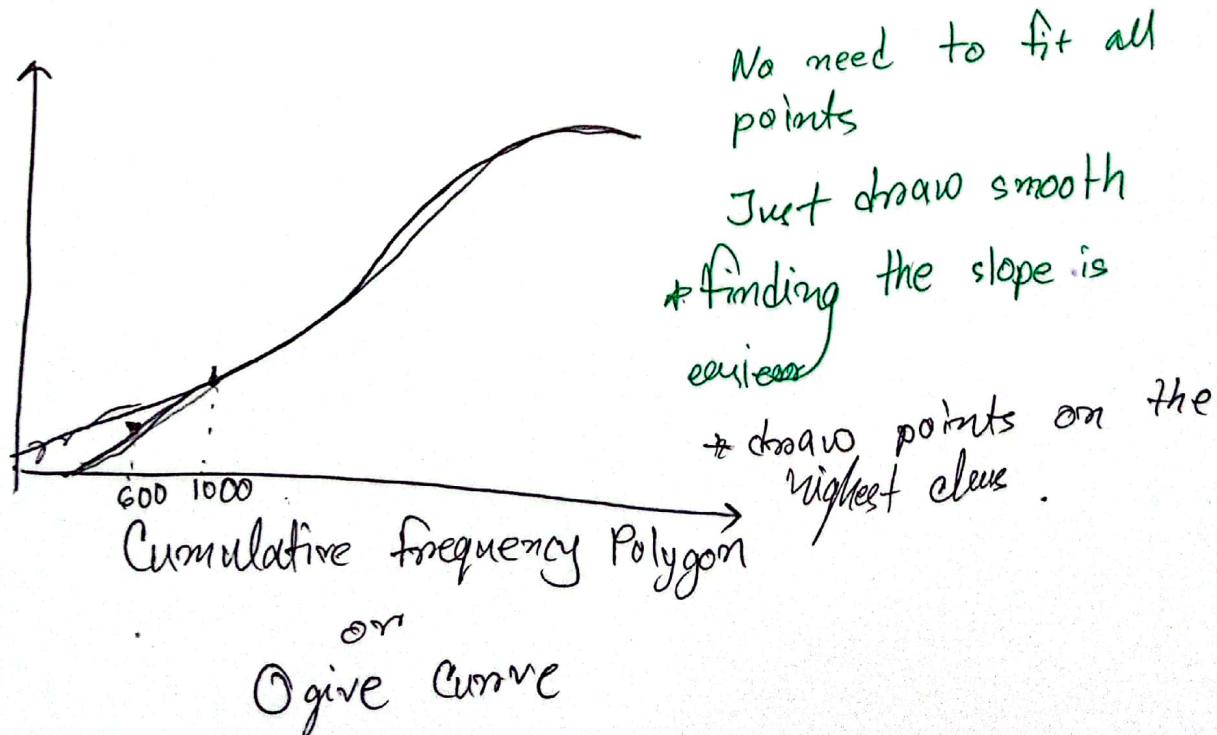


Fig: Histogram of the Profit -
(Table Name but replace it ~~for~~ with Histogram)

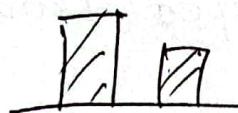


* When * clauses have names, we use bar diagrams or pie chart

Bar Diagram:

→ X-axis: Name ; scaling not needed

Y-axis: Frequency ; scaling needed

* Mark the bars : 

* Clauses without names can also be plotted.

Replace names with class width.

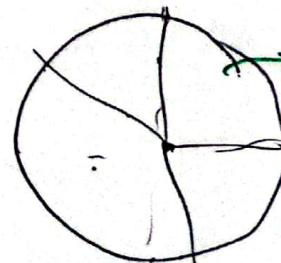
Y-axis remains same

Pie Chart:

At first we calculate the angles

Class	Frequency	Angle
Total	n	

* No need for protractor
or maintaining serial
or showing calculation



→ use eye to calculate

Measures of Central Tendency

1. Mean → Arithmetic mean
 → Geometric means (%)
 → Harmonic ($\frac{1}{x+y}$ units) $\overrightarrow{\text{rate}}$

2. Median

3. Mode

Formula for Disgrouped Data

Mean :

Arithmetic mean: $\bar{X} = \frac{\sum x_i}{n}$

Geometric mean: $\sqrt[n]{x_0 x_1 x_2 \dots x_n} = \prod (x_i)$

(Lim Marshall - 72)

(S.P Gupta - 14: 99)

Harmonic mean = $\frac{n}{\sum \frac{1}{x_i}}$

* Median:

The value which takes the middle most position
in ordered arranged data set.

* Mode:

The value with maximum numbers of occurrences/
frequency in a data set is known as mode.

Formula For Grouped Data:

* Arithmetic Mean: $\bar{X} = \frac{\sum f_i x_i}{N}$

* Geometric Mean: Antilog $\frac{1}{N} \sum f_i \log x_i$; \log_{10} (NOT \ln)
Antilog $x = 10^x$

* Harmonic Mean: $\frac{N}{\sum \frac{f_i}{x_i}}$

* Median: $M_d = L + \frac{\frac{N}{2} - C.F.}{f_m} \times C$

L = lower limit of median class
 $C.F.$ = cumulative frequency
 f_m = median class frequency
 C = median class width
 N = total observation

* Mode: $M_o = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times C$

L = lower limit

Δ_1 = diff. between frequencies of modal and premodal class

Δ_2 = "

"

" " modal and post modal class

C = modal class width

Example of Geo. Mean.

As the data is percentage, we will use geometric mean.

Let, replace with unit mentioned
in the question.

100 ~~is~~ [rate] is invested

return 30 %	meant 130 [made]
" 20 %	" 120 [rate]
" -40 %	" 80 [rate]
" 200 %	" 300 [rate]

i.e convert the data in smaller form.

That is: 1.3, 1.2, 0.6, 3.0

Then, the (geometric) mean rate of return is:

$$G.M = \sqrt[n]{(x_1)(x_2) \dots (x_n)} = \sqrt[4]{(1.3)(1.2)(0.6)(3.0)}$$

$$= 1.294$$

i.e. ~~129.4%~~

∴ Return is 29.4 %.

Example (S. G. Gupta) :

% Rise	Expenses taking 100%.	log X
32	132	
40	140	
:		

* Do : S. G. Gupta: Miscellaneous Examples

(27, 30, 32, 33)
(But do all)

Harmonic Mean:

$$H.M = \frac{1}{\frac{1}{x_1} + \frac{1}{x_2} + \frac{1}{x_3} + \dots}$$

$$= \frac{n}{\sum \frac{1}{x_i}}$$

Used whenever there is a per-unit ($/ \text{unit}$) involved

* Any rate (km/hr , takey/kg)

x goes 20 km per in 5 hours

$$\frac{20 \text{ km}}{5 \text{ hours}} = 4 \text{ km/hr} ; \left[\text{for finding distance in time} \right]$$

$$\frac{5 \text{ hours}}{2 \text{ km}} = \frac{1}{4} \text{ hr/km} ; \left[\text{for finding time required to cover } \text{distance} \right]$$

Example:

$$A \rightarrow 4 \text{ min/work} = 0.25$$

$$B \rightarrow 5 \text{ min/work} = 0.2$$

$$C \rightarrow 6 \text{ min/work} = 0.167$$

$$D \rightarrow 10 \text{ min/work} = 0.1$$

$$E \rightarrow 12 \text{ min/work} = 0.833$$

i) Avg work per min = ?

ii) How much work in 6 hours per day

The work is done per unit, we will use harmonic mean

$$\therefore H.M = \frac{5}{\frac{1}{4} + \frac{1}{5} + \frac{1}{6} + \frac{1}{10} + \frac{1}{12}} = 0.25 \text{ min/work}$$
$$= 0.16 \text{ work/min (Ans)}$$

$$\therefore \text{work in 6 hours} = 0.16 \times 80 \times 60$$

= .

Combined Mean

$$\begin{array}{|c|} \hline n_1 \\ \hline \bar{x}_1 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline n_2 \\ \hline \bar{x}_2 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline n \\ \hline \bar{x} \\ \hline \end{array}$$

k number of groups containing n_i elements with

\bar{x}_i mean

~~They will have form~~

∴ Combined Mean, $\bar{x} = \frac{\sum n_i \bar{x}_i}{\sum n_i}$

Example (6(a)) :

Given that,

$$n_1 = 100$$

$$n_2 = 80$$

$$\bar{x}_1 = 4570 \$$$

$$\bar{x}_2 = 750 \$$$

$$\therefore \bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2} = -$$

Given: 6(b)

$$\bar{x}_c = 72$$

$$n_1 + n_2 = 100$$

$$n_1 = 70$$

$$n_2 = 30$$

$$\left. \begin{array}{l} \bar{x}_1 = \text{something} \\ \bar{x}_2 = ? \end{array} \right\}$$

$$\bar{x}_c = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$

$$\therefore \bar{x}_2 = . . .$$

Q

Weighted Mean:

$$\bar{x}_w = \frac{\sum w_i x_i}{\sum w_i}$$

* Arithmetic Mean but FANCY ~

Q

Grouped Data

Example: Random sample of 25 from 2 brands
 Decide which brand performs better
in average

Hour	1-2	2-3	3-4	4-5	5-6
Type - I	2	4	8	7	4
Type - II	3	6	8	6	2

* Anything that has less deviation is more consistent
 ∴ DO NOT use this method.
 for arith mean ↑ for geo mean ↑ Harmonic mean ↑ median ↑
 ✓

Class	Frequency (f.)	x_i	$f_i x_i$	$f_i \log x_i$	$\frac{f_i}{\sum f_i}$	Cumulative frequency c.f.
Total	$N =$		$\sum f_i x_i$	$\sum f_i \log x_i$	$\sum \frac{f_i}{\sum f_i}$	

No approximation is allowed

BUT keep upto 3 digits
 i.e. $1.999999 \approx 2.000$

* ~~Avg.~~ ~~ALWAYS,~~

$$A.M > G.M > H.M$$

* Median class is the class with half of the total frequency

Consumption of fuel (lit/hour)	Number of generators (f _i)	x _i	f _i x _i	f _i log x _i	$\frac{f_i}{x_i}$	Cumulative Frequency c _i f _i
1 - 2	2	1.5	3	0.352	1.33	2
2 - 3	4	2.5	10	1.592	1.6	6
3 - 4	8	3.5	28	4.352	2.285	14
4 - 5	7	4.5	31.5	4.572	1.555	21
5 - 6	4	5.5	22	2.961	0.727	25
Total	N = 25		$\sum f_i x_i = 94.5$	$\sum f_i \log x_i = 13.829$	$\sum \frac{f_i}{x_i} = 7.492$	

$$\therefore \text{Arithmetic mean } \bar{x} = \frac{\sum f_i x_i}{N} = \frac{94.5}{25}$$

$$= 3.78 \text{ lit/hour}$$

$$\therefore \text{Geo. Mean} = \text{Antilog} \left\{ \frac{1}{N} \sum f_i \log x_i \right\}$$

$$= \text{Antilog} \left\{ \frac{1}{25} \times 13.829 \right\}$$

$$= 3.574 \text{ lit/hour}$$

Harmonic mean, $H.M = \frac{N}{\sum \frac{f_i}{x_i}} = \frac{25}{7.497}$
 $= 3.3347 \text{ lit/hm.}$

Hence, we have $N = 25$ observations and $N/2$
 $= 25/2 = 12.5$

Thus, from the table we may decide median class
 is 3-4.

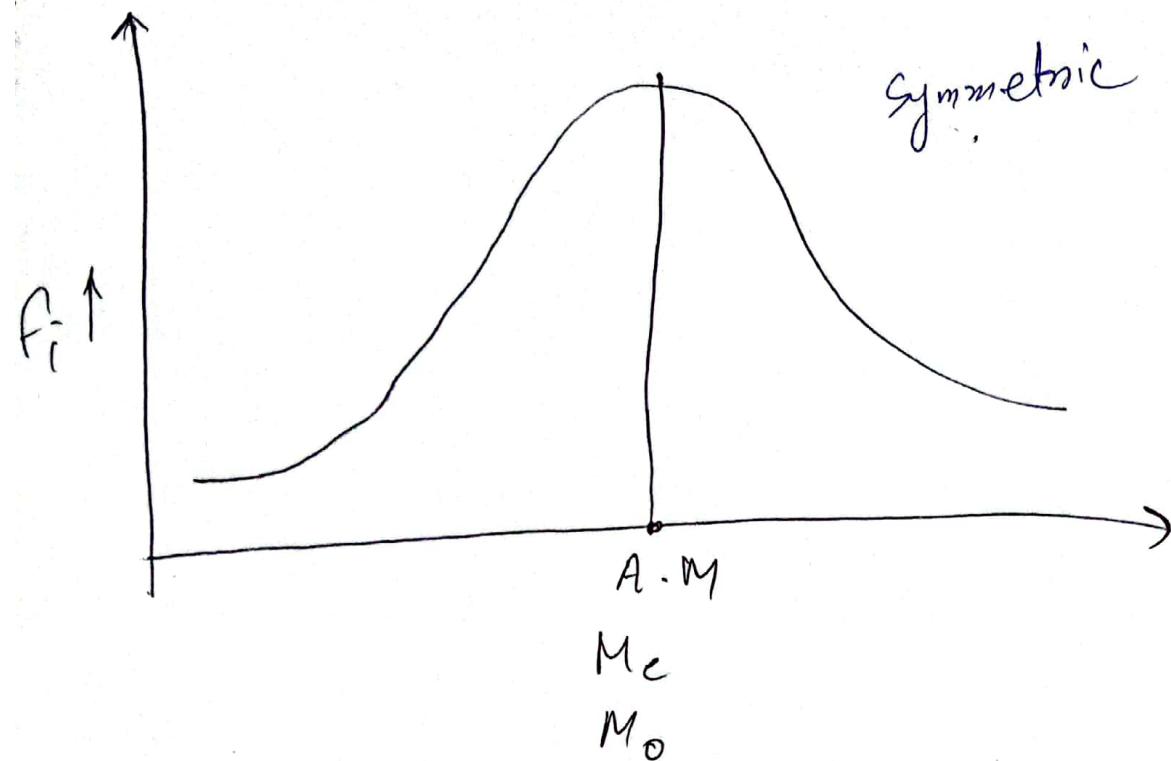
Therefore, the median is, $M_d = L + \frac{\frac{N}{2} - C.P}{f_m} \times C$
 $= 3 + \frac{12.5 - 6}{8} \times 1 =$

$$= 3.08125 \text{ lit/hm}$$

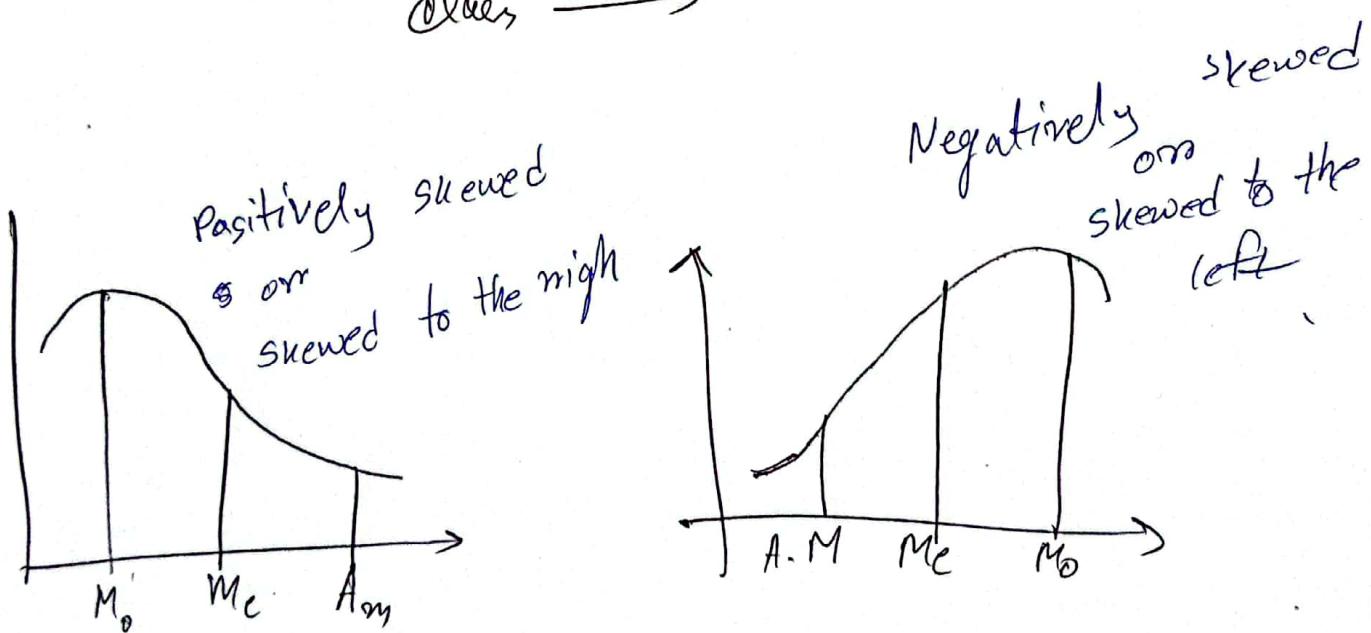
Hence, the class with the maximum number of frequencies
 is 3-4, which is 8. Thus, 3-4

is own modal class

The mode is, $M_o = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times C$
 $= 3 + \frac{8-4}{(8-4)+(8-2)} \times 1 = 3 + 0.8$
 $= 3.8 \text{ lit/hm.}$



class →



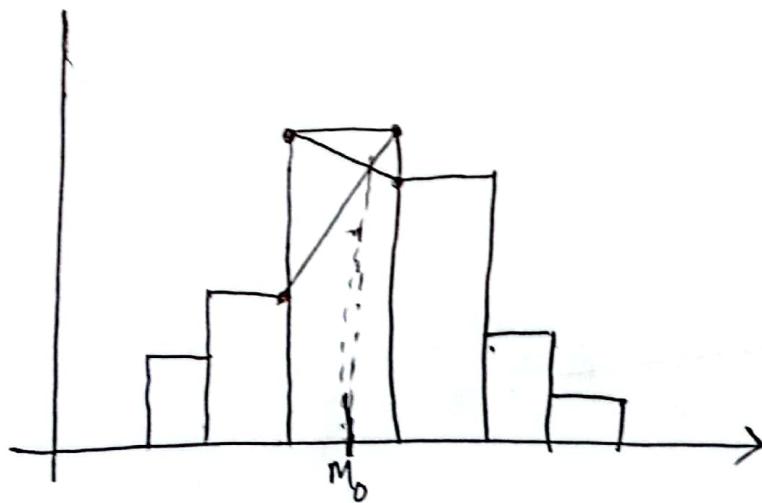
Relative Positions of the Mean, Median and Mode

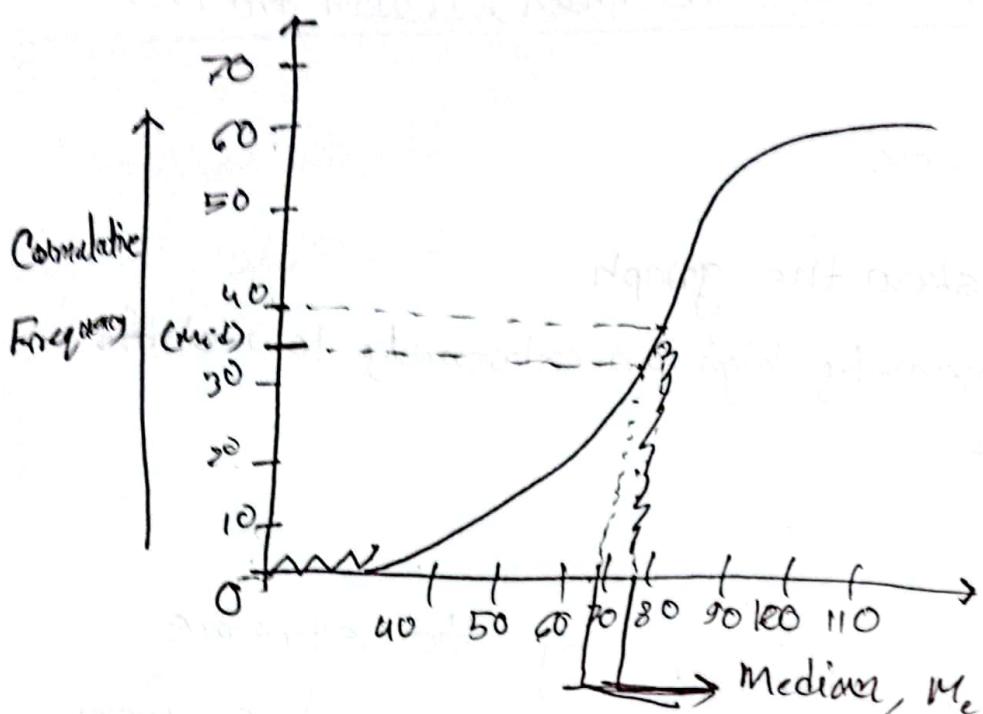
Box Removal of Skewness

- * Extreme values skew the graph
- * Values that are extremely high or extremely low shifts the value towards it.

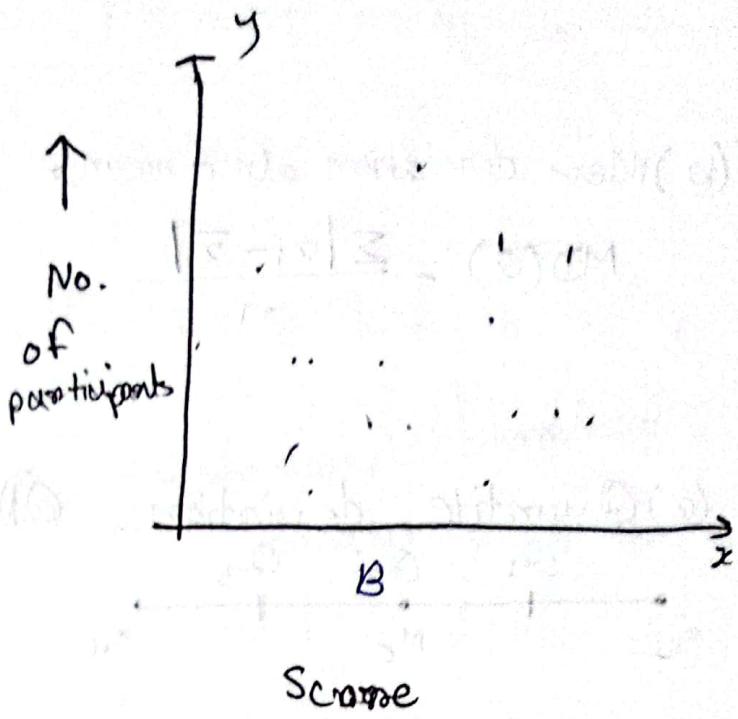
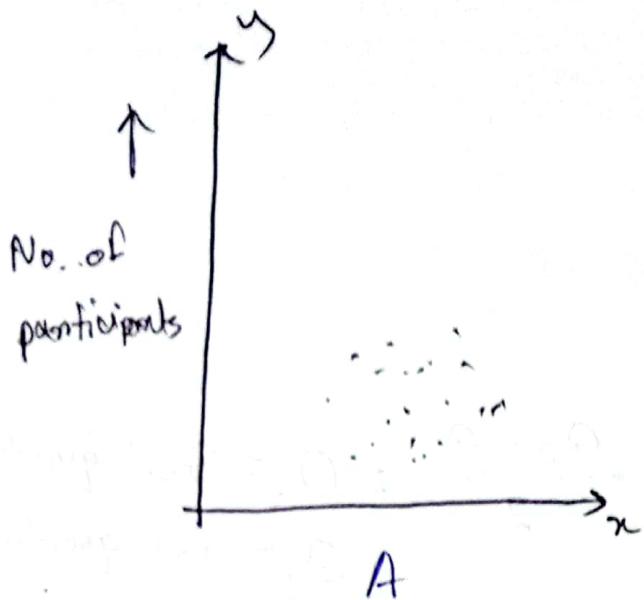
*+

- * Arithmetic mean is affected by the extreme value in a data set while the median is free from it (- Explain with example)





Measures of Dispersion:



A is more consistent/reliable/steady etc than B

Two types of measures:

- i) Absolute measures
- ii) Relative "

- (a) Range (R)
- (b) Mean deviation about mean ($M.D., \bar{x}$)
- (c) Quartile deviation (Q.D.)
- (d) Standard " (S)
- (a) Coefficient of range (C.R.)
- (b) " " mean deviation about mean ($C.M.D., \bar{x}$)
- (c) Coefficient of Quartile deviation (C.Q.D.)
- (d) Coefficient of standard deviation
- or, Coefficient of Variation (C.V.)

Two long coupled:

Range = $R = x_H - x_L$; x_H = highest value
 x_L = lowest value

~~*(b)~~ Mean deviation about mean

$$MD(\bar{x}) = \frac{\sum |x_i - \bar{x}|}{n}$$

(e) Quantile deviation, $QD = \frac{Q_3 - Q_1}{2}$; Q_3 = 3rd quartile
 Q_1 = 1st quartile





(d) Standard Deviation, $s = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n-1}}$; $\bar{x} = \frac{\sum x}{n}$

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

↳ preferred

Fors goduped:

(i) (a) No need

$$*(b) MD(\bar{x}) = \frac{\sum f_i |x_i - \bar{x}|}{N}; \quad N = \sum f_i \\ \bar{x} = \frac{\sum f_i x_i}{N}$$

(c) No need

$$\text{Ans} \quad (d) S = \sqrt{\frac{\sum f_i (x_i - \bar{x})^2}{N-1}} = \sqrt{\frac{\sum f_i x_i^2 - (\sum f_i x_i)^2 / N}{N-1}}$$

↳ preferred

$$2(i)(a) C.R. = \frac{x_H - x_L}{x_H + x_L} \times 100$$

$$(b) C.MD(\bar{x}) = \frac{MD(\bar{x})}{\bar{x}} \times 100$$

$$(c) CQD = \frac{Q_3 - Q_1}{Q_3 + Q_1} \times 100$$

$$(d) C.V = \frac{s}{\bar{x}} \times 100; s = \text{sample standard deviation}$$

23.5	23.8	24.2	24.5	24.8	25.2	25.5	25.8	26.2	26.5	26.8	27.2	27.5	27.8	28.2	28.5	28.8	29.2	29.5	29.8	30.2	30.5	30.8	31.2	31.5	31.8	32.2	32.5	32.8	33.2	33.5	33.8	34.2	34.5	34.8	35.2	35.5	35.8	36.2	36.5	36.8	37.2	37.5	37.8	38.2	38.5	38.8	39.2	39.5	39.8	40.2	40.5	40.8	41.2	41.5	41.8	42.2	42.5	42.8	43.2	43.5	43.8	44.2	44.5	44.8	45.2	45.5	45.8	46.2	46.5	46.8	47.2	47.5	47.8	48.2	48.5	48.8	49.2	49.5	49.8	50.2	50.5	50.8	51.2	51.5	51.8	52.2	52.5	52.8	53.2	53.5	53.8	54.2	54.5	54.8	55.2	55.5	55.8	56.2	56.5	56.8	57.2	57.5	57.8	58.2	58.5	58.8	59.2	59.5	59.8	60.2	60.5	60.8	61.2	61.5	61.8	62.2	62.5	62.8	63.2	63.5	63.8	64.2	64.5	64.8	65.2	65.5	65.8	66.2	66.5	66.8	67.2	67.5	67.8	68.2	68.5	68.8	69.2	69.5	69.8	70.2	70.5	70.8	71.2	71.5	71.8	72.2	72.5	72.8	73.2	73.5	73.8	74.2	74.5	74.8	75.2	75.5	75.8	76.2	76.5	76.8	77.2	77.5	77.8	78.2	78.5	78.8	79.2	79.5	79.8	80.2	80.5	80.8	81.2	81.5	81.8	82.2	82.5	82.8	83.2	83.5	83.8	84.2	84.5	84.8	85.2	85.5	85.8	86.2	86.5	86.8	87.2	87.5	87.8	88.2	88.5	88.8	89.2	89.5	89.8	90.2	90.5	90.8	91.2	91.5	91.8	92.2	92.5	92.8	93.2	93.5	93.8	94.2	94.5	94.8	95.2	95.5	95.8	96.2	96.5	96.8	97.2	97.5	97.8	98.2	98.5	98.8	99.2	99.5	99.8	100.2	100.5	100.8	101.2	101.5	101.8	102.2	102.5	102.8	103.2	103.5	103.8	104.2	104.5	104.8	105.2	105.5	105.8	106.2	106.5	106.8	107.2	107.5	107.8	108.2	108.5	108.8	109.2	109.5	109.8	110.2	110.5	110.8	111.2	111.5	111.8	112.2	112.5	112.8	113.2	113.5	113.8	114.2	114.5	114.8	115.2	115.5	115.8	116.2	116.5	116.8	117.2	117.5	117.8	118.2	118.5	118.8	119.2	119.5	119.8	120.2	120.5	120.8	121.2	121.5	121.8	122.2	122.5	122.8	123.2	123.5	123.8	124.2	124.5	124.8	125.2	125.5	125.8	126.2	126.5	126.8	127.2	127.5	127.8	128.2	128.5	128.8	129.2	129.5	129.8	130.
------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	------

#Example: Suppose, the value from strength test on some samples of yarns from different spinning machines are given (in lbs.)

from 1st machine:

8.5, 8.5, 8.8, 7.8, 7.9, 9.15, 9.0

from 2nd machine:

8.8, 9.2, 8.4, 8.4, 9.2, 9.10, 8.25

Decide which machine looks more reliable regarding their performance variations.

Solⁿ: The calculation table for standard deviation is given below (for first machine)

x_i	x_i^2
8.5	72.25
8.5	72.25
8.8	77.44
7.8	60.84
7.9	62.41
9.15	83.7225
9.0	81
$\sum x_i = 51.15$	$\sum x_i^2 = 509.9125$

$$\sum x_i = 51.15 \quad \sum x_i^2 = 509.9125$$

∴ the standard deviation for the first spinning machine

$$\text{is, } S = \sqrt{\frac{\sum x_i^2 - (\sum x_i)^2}{n-1}} = 4.7636$$

That is, the strengths of yarns are varying by ~~17%~~ 4.763% from the value of the overall outputs.

Now, the coefficient of variation is,

$$C.V = \frac{\sigma}{\bar{x}} \times 100$$

$$= \boxed{20} 65\%$$

$$= \cancel{65\%} 65\%$$

That is, the percentage rate of variation of strengths of yarns is 12% for the 1st machine.

Do for 2nd machine]

The calculation table for standard deviation is given below

(for second machine):

x_i	x_i^2
8.8	77.44
9.2	84.64
8.4	70.56
8.4	70.56
9.2	84.64
9.10	82.81
8.25	68.0625
$\sum x_i = 61.35$	
$\sum x_i^2 = 538.725$	

The standard deviation for the second machine is,

$$S = \sqrt{\frac{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}{n-1}}$$

$$= \sqrt{\frac{(538.7125) - \frac{(61.35)^2}{7}}{6}}$$

$$= 0.413$$

$$C.V_2 = \frac{s}{\bar{x}} \times 100 \text{ \%}$$

$$= \frac{0.413}{8.764} \times 100 \text{ \%}$$

$$= 4.7125 \text{ \%}$$

$$\therefore C.V_2 < C.V_1$$

∴ Machine 2nd Machine looks more reliable.

C-9(W-3) *Gupta - Ch-5 (144) 30/5/24
 (#Slide Example) (148) → avoid combined deviation (online)
 (33, 37)

	65-70	70-75	75-85	85-95	95-100
Machine I	1	4	8	6	2
Machine II	2	6	8	4	1

Calculation table for machine - I

Class interval	f _i (Frequency)	x _i	f _i x _i	f _i x _i ²
65-70	1	67.5	67.5	4556.25
70-75	4	72.5	290	21025
75-85	8	80	640	51200
85-95	6	90	540	
			1732.5	
N = 21			$\sum f_i x_i = 1732.5$	$\sum f_i x_i^2 = 144393.75$

$$S = \sqrt{\frac{\sum f_i x_i^2 - \frac{(\sum f_i x_i)^2}{N}}{N-1}} = 8.5513 \text{ percent spot removal}$$

which means, the amount of variation in spot removal perc is 8.5513%

$$\therefore \bar{x} = \frac{\sum f_i x_i}{N} = \frac{1732.5}{21} = 82.5$$

$$\therefore CV = \frac{S}{\bar{x}} \times 100 \% = 10.36 \%$$

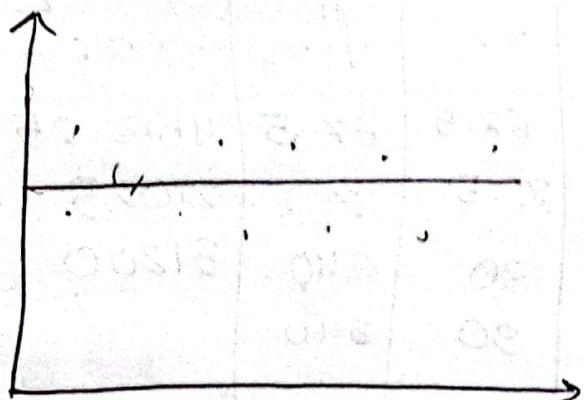
That is, the percentage rate of variation in the spot removals of machine-I is 10.36%
 [Do point 2]

Simple Linear Regression and Correlation

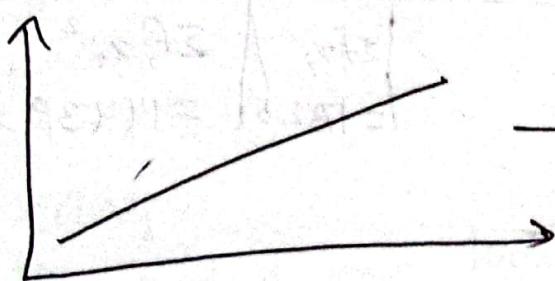
Correlation: relation between two or more variables

- a) Simple Correlation: Correlation between only two variables
- b) Multiple Correlation: Correlation " more than " "
- c) Partial: Ignore what you don't want

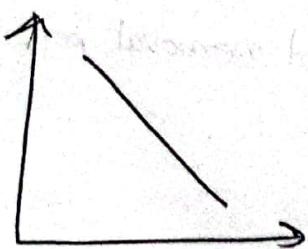
* Graphs with dots are called
Scatter plots



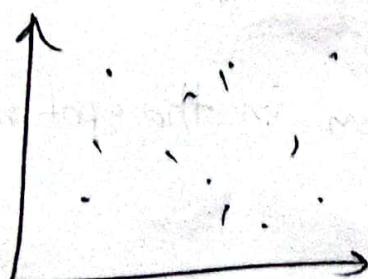
→ zero correlation



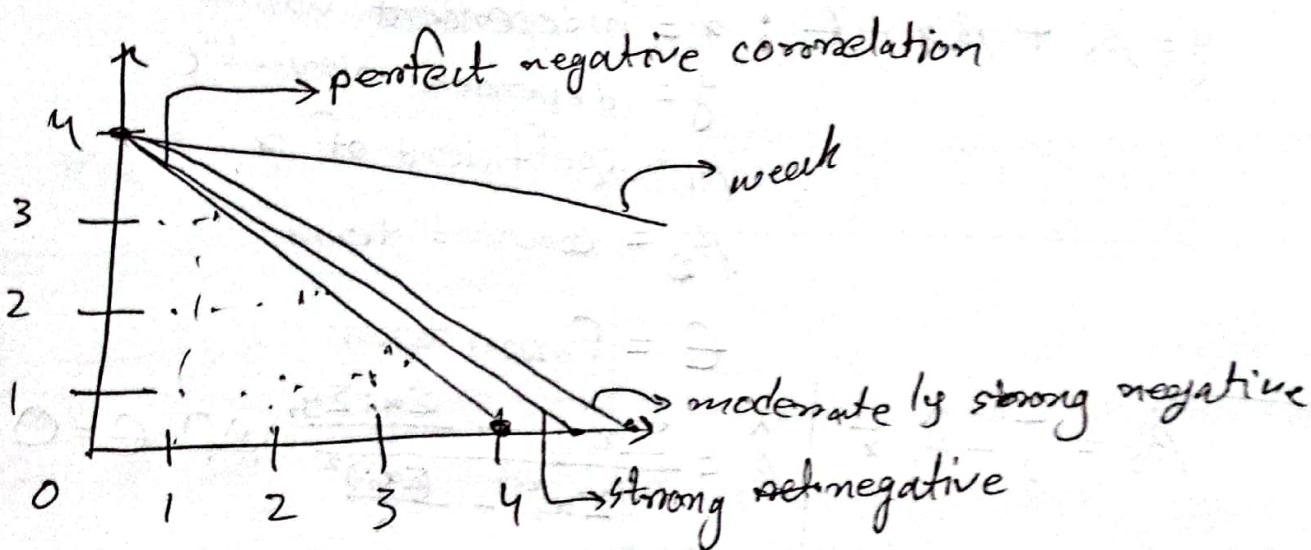
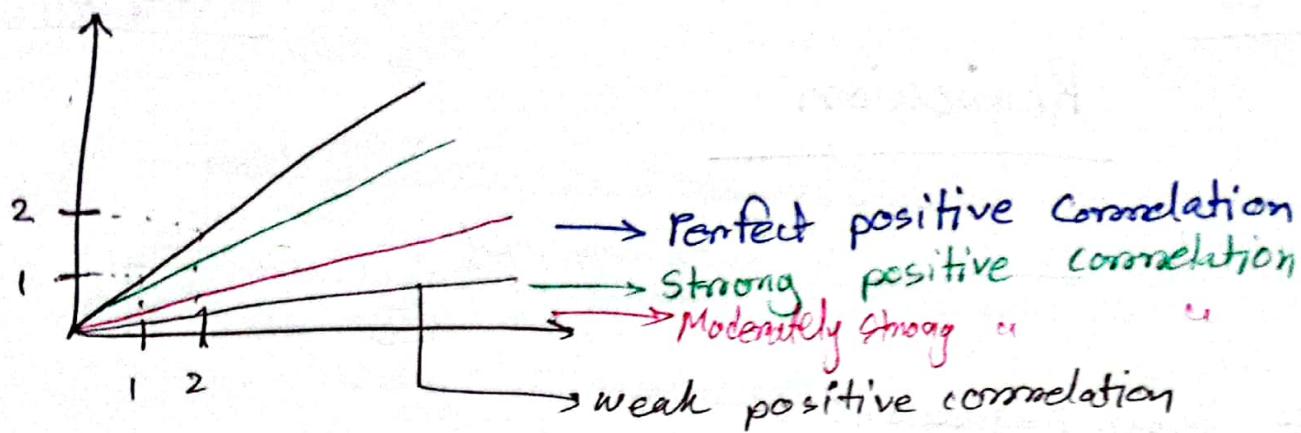
→ positive correlation



Negative correlation



→ No significant evidence of correlation



Pearson's Correlation Coefficient,

$$r = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sqrt{\left\{ \sum x_i^2 - \frac{(\sum x_i)^2}{n} \right\} \left\{ \sum y_i^2 - \frac{(\sum y_i)^2}{n} \right\}}}$$

Covariance (x, y)
 determines
 +ve or -ve

perfect $\rightarrow 1$
 $0.7 < 1 \rightarrow$ strong positive
 $0.4 < 0.7 \rightarrow$ moderately strong
 $> 0 < 0.4 \rightarrow$ weak

$-1 \leq r \leq 1$; $r = 0$

Regression

* If indicates _____ is working well

Simple linear regression:

$$y = \beta_0 + \beta_1 x + \epsilon ; x = \text{independent variable}$$

$y = \text{dependent variable}$

$\beta_1 = \text{coefficient of } x$



$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \left| \quad \hat{\beta}_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} \right. \cdot \text{If } \epsilon = 0$$

* Error will happen naturally

Regression:

Let,

$$X = [\text{data}] \text{ (unit)}$$

$$Y = [\text{data}] \text{ (unit)}$$

The simple linear regression equation is, → The formula from last page

The calculation table of regression model/parameters is given below:

x_i	y_i	x_i^2	y_i^2	$x_i y_i$
$\sum x_i =$	$\sum y_i =$	$\sum x_i^2 =$	$\sum y_i^2 =$	$\sum x_i y_i =$
10	10	100	100	100
20	20	400	400	400
30	30	900	900	900
40	40	1600	1600	1600
50	50	2500	2500	2500
60	60	3600	3600	3600
70	70	4900	4900	4900
80	80	6400	6400	6400
90	90	8100	8100	8100
100	100	10000	10000	10000

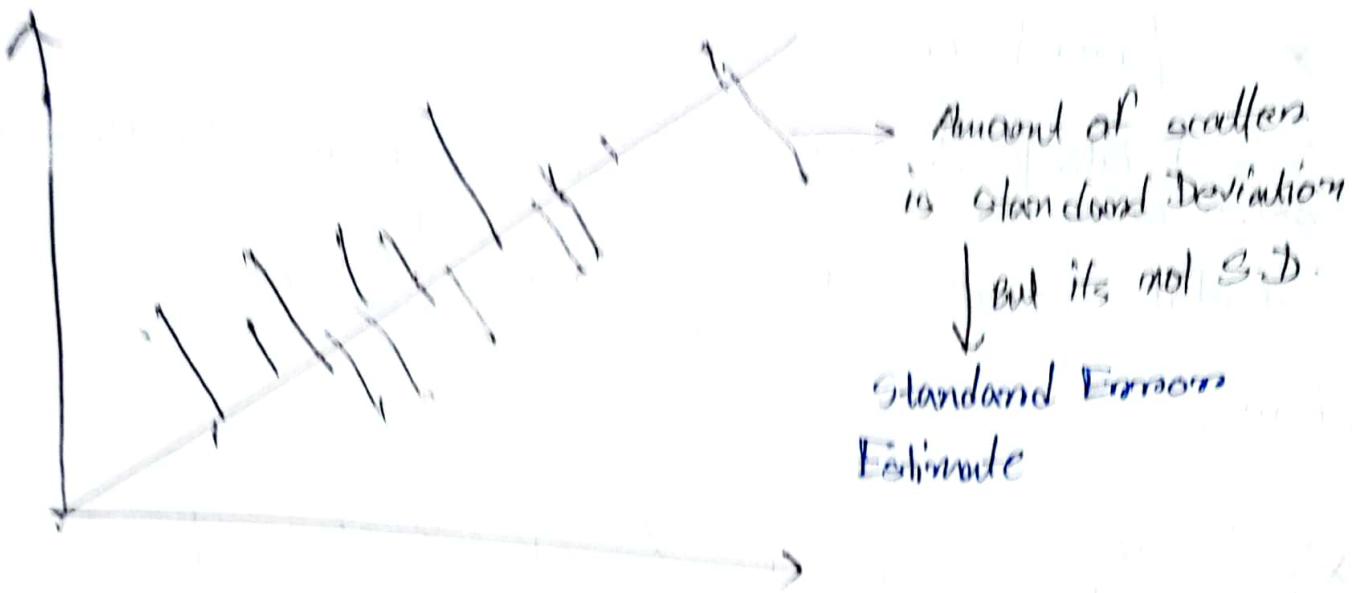
Thus, the estimates of the regression parameters are,

$$\hat{\beta}_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = []$$

$$\hat{\beta}_0 = []$$

The regression model is,

$$y = -33 + 0.7x$$



Standard Error of the Estimate:

The measure of the dispersion or scatter of the observed values (sample observations) around the estimated regression line.

$$S_{yx} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} = \sqrt{\frac{\sum y_i^2 - \hat{\beta}_0 \sum y_i - \hat{\beta}_1 \sum x_i y_i}{n-2}}$$

↑ preferred

$$= \sqrt{MSE} = \sqrt{\text{Mean Square Error}}$$

Coefficient of Determination

The amount of variation in the response variable (y) that can be explained by the variation of independent variable (x).

$$R^2 = \frac{RSS}{TSS} = 1 - \frac{ESS}{TSS};$$

RSS = Regression sum of square

$$\text{square} = \sum (\hat{y}_i - \bar{y})^2$$

ESS = Error sum of square

$$\text{square} = \sum (y_i - \hat{y}_i)^2 = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$$

Here, $TSS = RSS + ESS$ [OSR \Rightarrow preferred]

where, TSS = Total sum of square = $\sum (y_i - \bar{y})^2 = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$

Probability

* Ask P

* Randomness = uncertain outcome

* Variables used in probability are called Random Variables

* Sample Space : Set of all possible outcomes. $\rightarrow S$

* Event : Subset of sample space ~~for from~~ that contains required outcomes.

* Discrete : Individual

Continuous : Between

* Mutually Exclusive : Two or more events with no common event
 $\hookrightarrow A \cap B = \emptyset$

Probability :

(i) Let, a random experiment has 'n' outcomes
 A be = any specific characteristic
 n_A = number of outcomes having characteristic 'A'
 $P(A) = \frac{n_A}{n}$



Axioms of Probability:

Let,

a random experiment has 'n' outcomes:

$$A_1, A_2, A_3, \dots, A_n$$

(i) $\sum P(A_i) = 1$

(ii) $0 \leq P(A_i) \leq 1$ for any A_i

(iii) For any two events A_1 and A_2 , if $A_1 \cap A_2 \neq \emptyset$, then,

$$P(A_1 \cup A_2) = P(A_1) + P(A_2)$$

Rules of Solving Problems:

(i) Complementary Rule: $P(A) = 1 - P(A')$

(ii) Additional Rule: For any two events A and B ,

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ &= P(A) + P(B); \text{ [mutually exclusive]} \end{aligned}$$

(iii) Conditional Rule: If A depends on B then,

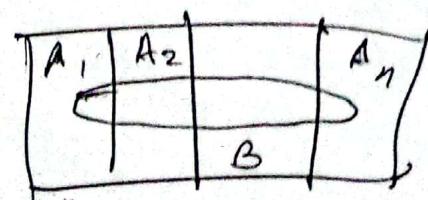
$$A|B \equiv A \text{ given } B \text{ depends on } P(A|B) = \frac{P(A \cap B)}{P(B)}; [P(B) \neq 0]$$

(iv) Independence Rule: Any two events A and B are said to be independent of each other if,

$$P(A \cap B) = P(A) \cdot P(B)$$

(v) Bayes' Rule: $P(A_i|B) = \frac{P(A_i) \cdot P(B|A_i)}{\sum P(A_i)P(B|A_i)}$

↳ total sum



Suppose, in a controlled room, 3 operators are giving commands for completing a mission. Their previous records say that:

The chances of failures of the three operators are 9.1%, 12% and 7.5%.

The given mission will be considered completed if the 3 ops work together and each operator is independent of the other two ops.

Find the probability that a newly given mission will be completed successfully by them.

Let,

A = first operator works successfully

" "

B = second "

" "

C = third "

Given that,

their failure rates in terms of probability are:

$$P(A') = 0.09, P(B') = 0.12, P(C') = 0.075$$

$$\begin{aligned} \text{Their success rates} &= P(A) + P(B) + P(C) \\ &= \{1 - P(A')\} \{1 - P(B')\} \{1 - P(C')\} \\ &= 1 - 0.09 = 1 - 0.12 = 1 - 0.075 \\ &= 0.91 = 0.88 = 0.925 \end{aligned}$$

The probability that the given mission will be completed successfully is,

$$\begin{aligned} P(A \cap B \cap C) &= P(A) \cdot P(B) \cdot P(C) \\ &= (0.91)(0.82)(0.925); [] \\ &= 0.74074 \\ &= 74.074\% \quad (\text{Ans.}) \end{aligned}$$

+ 25/09/22

\cup { either or, neither nor,
minimum of, at least }

\cap { and, all, both, including all,
all inclusive, nothing will be left alone etc. }

such that, given that, subjected to the,
conditional condition that, restricted to, constrained to etc
Rule

C-21 (w-14)

25/09/24

Actual data

Contamination	Center	Edge	Total	(Ch-2, Ex-2.19)
Low	514	68	582	Suppose Si wafers
High	112	246	358	= area to setup.
Total	626	314	940	In total 582 area of low level contamination, 112 wafers are high but center. These are Random. 84314 edge ones.

Suppose one wafer is selected at

Find the probability that

- i) It is located in the edge area or it has low level contamination
- (ii) It is located in the edge area and it has low level contamination
- (iii) It is located in the edge area such that it has low level contamination but it has
- (iv) It is located in the edge area but it has low level contamination

The inform

$$(103)^2 = (107 + 3)^2 = 107^2 + 2 \cdot 107 \cdot 3 + 3^2 = 107^2 + 62 + 9 = 107^2 + 63$$

Sol^m: The int Let,

H = high level contamination

L = low "

C = center area

E = edge

The information table is,

Contamination Table	location			Total
	Edge (E)	Center (C)	Total	
High (H)	286	112	358	
Low (L)	68	514	582	
Total	314	626	940	

Now,

the probability that a randomly selected wafer is

(i) located in the edge area or it has low level contamination

$$\begin{aligned} P(E \cup L) &= P(E) + P(L) - P(E \cap L) \\ &= \frac{314}{940} + \frac{582}{940} - \frac{68}{940} = \end{aligned}$$

(ii) Located in the edge and it has low level
contamination = $P(E \cap L) = \frac{68}{940} =$

(iii) Located in the edge such it has low level

$$\text{contamination} = P(E|L) = \frac{P(E \cap L)}{P(L)}$$
$$= \frac{\frac{68}{940}}{\frac{582}{940}}$$

$$= \frac{68}{582}$$

(iv) Same as (ii)

* Most something (Red book): 2.20, 2.22, 2.9

Exercise: 2.78, 2.77, 2.81,
2.90, 2.94

2. Z. Bayes Theorem:

A_1	A_2	A_3
35%.	40%	25%.
0.01% errors	0.011% errors	0.012% errors

Suppose there is an error.

Which person has the most probability of doing the error?

Solⁿ:

Let,

A_1 = The work is done by the first person

A_2 = " " " second "

A_3 = " " " third "

and E = The work is error

Given that,

the overall probability of work of the three persons
are

$$P(A_1) = 0.35$$

$$P(A_2) = 0.4$$

$$P(A_3) = 0.25$$

The errors rates of the three persons are:

$$P(E|A_1) = 0.0001$$

$$P(E|A_2) = 0.00011$$

$$P(E|A_3) = 0.00012$$

Then, the probability that a randomly found error in the work of a given day is done by first person is

$$P(A_1|E) = \frac{P(A_1) \times P(E|A_1)}{\sum P(A_i) P(E|A_i)}$$

$$= \frac{P_1(A_1) \times P(E|A_1)}{P(A_1)P(E|A_1) + P(A_2)P(E|A_2) + P(A_3)P(E|A_3)}$$

$$= \frac{0.35 \times 0.0001}{(0.35)(0.0001) + (0.4)(0.00011) + (0.25)(0.00012)}$$

$$2^{\text{nd}} \text{ Person} = \frac{P(A_2) P(E|A_2)}{\boxed{P(A_1) P(E|A_1) + P(A_2) P(E|A_2)}}$$

$$3^{\text{rd}} \text{ Person} = \frac{P(A_3) P(E|A_3)}{\boxed{}}$$

Maybe, Guaranteed value is most responsible for the errors.

2. 38, 2.148,

SP Gupta: Ch - 11: 400 \rightarrow 14, 15, 27, 30,
 Exercise \rightarrow 29, 30, 31

$$\frac{(GA)^9 \times (GA)^9}{(GA)^9 \times (GA)^9} = (GA)^9$$

$$(GA)^9 \times (A)^9$$

$$(GA)^9 (GA)^9 + (GA)^9 (A)^9$$

$$(GA)^9 (GA)^9 +$$

$$1000.0 \times 38.0$$

$$\frac{1000.0 \times 38.0 + 1000.0 \times 22.0}{1000.0 \times 22.0} =$$

Probability Distributions

Binomial Distribution

summation of Bernoulli Trial is Binomial Distribution

$$f(x) \rightarrow P(x) = \binom{n}{x} p^x q^{n-x}; x=0, 1, 2, \dots, n$$

n = trial numbers

n = No. of success

p = Probability of success (from past experience)

q = " failure

$$\text{and, } p+q=1$$

$$\text{Mean, } E(x) = np$$

$$\text{Variance, } \textcircled{2} \quad V(x) = npq$$

C-22(00-IU)

26/09/22

Bernoulli Trial:

Trial with only two possible outcomes; expected or failure.

Binomial Distribution:

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

where,

n = trial number

x = no. of success

p = probability of success (from prev. experience)

→ proba x is whatever you want to find

$q =$ " failure

and, ~~$p+q=1$~~

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

$$\text{Variance, } V(x) = npq \quad \left\{ \begin{array}{l} np > npq \end{array} \right.$$

The chance that a bit transmitted through a digital transmission channel is received in errors is 0.1. Also assume that the transmission trials are independent.

If 4 bits are transmitted, find the probabilities that

- Exactly 2 bits will be errors
- At least 2 bits will "
- At most 1 bits "
- More than 1 but not more than 3 bits will be errors.
- At best 3 bits will be correct

Solⁿ:

Let,

x be the number of bits received in errors
which follows a Binomial distribution

Given that,

trial number, $n = 4$

Probability of success, $p = 0.1$

" failure, $q = 1 - p = 0.9$

Thus,

Probability function is,
 $P(x) = \binom{n}{x} p^x q^{n-x} = \binom{4}{x} (0.1)^x (0.9)^{4-x}; [x=0, \dots, 4]$

(ii)

• Find the probabilities that,

(i) Exactly 2 bits will be in error

$$= P(2) = {}^4C_2 (0.1)^2 (0.9)^{4-2}$$

(ii) At least 2 bits will be in error

$$= P(x \geq 2) = P(2) + P(3) + P(4)$$

=

=

(iii) At most 1 bits in error

$$= P(x \leq 1) = P(0) + P(1)$$

=

$$= (0.9)^4 (1-0.9) \binom{4}{1} = 0.9^4 \cdot 0.1 \cdot 4 = 0.0324$$

(iv) More than 1 bit less than 3,

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

=

(v) Hence,

x is the number of bits received correctly

∴ Probability of success, $p = 0.9$

∴ " failure, $q = 0.1$

$$\therefore P(x) = {}^4C_x (0.9)^x (0.1)^{4-x} \quad [x = 0, 1, 2, 3, 4]$$

Now,

Probability of at least 3 bits will be received

$$\text{correctly} = P(x \leq 3) = P(0) + P(1) + P(2) + P(3)$$

$$= 1 - P(4) =$$

Mont $\rightarrow 3.18$,

$$Ex \rightarrow 3.84, 3.87, 3.91$$

$$SPGupta \rightarrow 421 \text{ (Ch - 12)} \rightarrow 4, 5, 35$$

Normal Distribution

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

Here,

Mean, $E(x) = \mu$; $[-\infty < \mu < \infty]$

Variance, $V(x) = \sigma^2$; $0 < \sigma^2 < \infty$

Standard deviation = σ

Standard Normal Distribution:

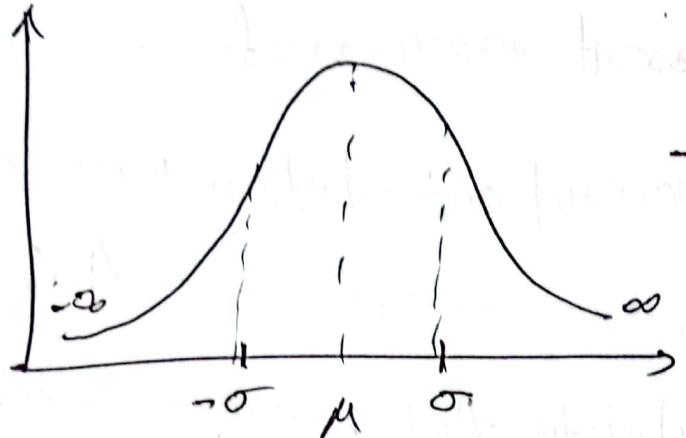
When $\mu=0$ and $\sigma^2=1$, then

then, it is a special case of Normal Distribution,

which is called standard normal distribution.

The variable is z where, $z = \frac{x-\mu}{\sigma}$

$$= (0) - 1 =$$



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

$$\int_{\mu-\sigma}^{\mu+\sigma} f(x) dx = 0.6827 ; \mu \pm \sigma \text{ covers } 68.27\% \text{ of area}$$

$$\int_{\mu-2\sigma}^{\mu+2\sigma} f(x) dx = 0.9545 ; \mu \pm 2\sigma \quad 95.45\% \text{ " "}$$

$$\int_{\mu-3\sigma}^{\mu+3\sigma} f(x) dx = 0.9972 ; \mu \pm 3\sigma \quad 99.72\% \text{ " "}$$

4.10, 4.13, 4.14

Assume that the current measurements in a strip of wire follow a normal distribution with a mean of 10 milliamperes and variance 4 (mA^2)

- (i) what is the probability that a measurement exceeds 13 mA?
- (ii) what is the probability that a current measurement is between 9 and 11 mA?
- (iii) Determine the value for which the probability that a current measurement is below this value is 0.98

[Use the stat table]

Mark — 4.15, 4.16
Ex → 4.58, 4.61, 4.62, 4.64, 4.65

Soln:

~~Q~~

Let,

x be the current measurement
which follows normal distribution

Given that,

mean, $\mu = 10 \text{ mA}$

Variance, $\sigma^2 = 4 \text{ mA}^2$

\therefore Standard deviation, $\sigma = 2 \text{ mA}$

The probability that the current measurement for
given a randomly given strip of wire is

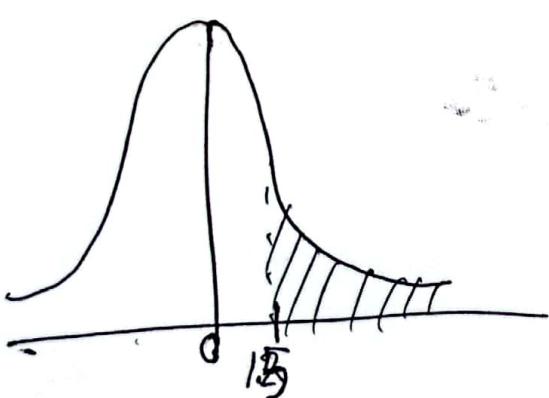
(i) will exceed 13 is,

$$P(x > 13) = P\left(\frac{x-\mu}{\sigma} > \frac{13-\mu}{\sigma}\right)$$

$$= P(z) = P(z > 1.5) = 1 - P(z < 1.5)$$

$$= 1 - 0.9332$$

$$= 0.0668$$



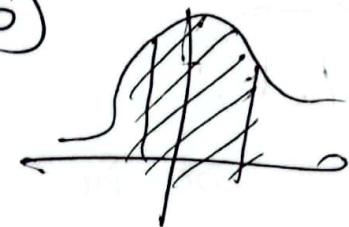
(ii) between 9 and 11

$$P(9 < z < 11) = P(-0.5 < z < 0.5)$$

$$= P(z < 0.5) - P(z < -0.5)$$

$$= 0.6915 - 0.3085$$

$$= 0.3830$$



(iii) the value for which the probability that a current is below the value is 0.98

$$z = 2.05$$

$$\text{or, } \frac{x-\mu}{\sigma} = 2.05$$

$$\therefore x = 14.1 \text{ mA}$$

$$(0.1 > 0.98) \Rightarrow 0.02 = (0.1 < 0.98) = 0.02$$

$$1886.0 - 1 =$$

$$1885.0 =$$

